

TRELLIS-CODED MODULATION ON TIME-DISPERSIVE CHANNELS

**A thesis presented for the degree of
Doctor of Philosophy
in
Electrical and Electronic Engineering
at the
University of Canterbury,
Christchurch,
New Zealand.**

**by
Chris J. Carlisle
B.E. (Hons 1)
December 1990**

Abstract

In this thesis we examine the performance of trellis-coded modulation (TCM) on time-dispersive channels. More specifically, we are interested in the performance improvements that TCM can offer in the presence of the residual intersymbol interference (ISI) that remains after non-ideal equalization on digital microwave radio (DMR) channels. The results are, however, applicable to other time-dispersive channels.

The performance of TCM on additive white Gaussian noise (AWGN) channels is well understood, and tight analytical bounds exist on the probability of the Viterbi decoder making a decision error. When a channel is also time-dispersive, the performance of TCM systems has, in the past, been studied mainly by simulation, due to the difficulty of formulating tractable analytical bounds on the error probability. In the work reported here, both simulation and analytical techniques are used.

The results of the simulation study show that TCM can improve the performance of a system with residual ISI. Although significant coding gains are achieved, the improvements in link outage are small, but useful. Simulation, however, is limited to symbol error probabilities greater than about 10^{-5} , and is not a particularly useful tool for estimating error probabilities over the range required for designing codes. There is a need for tight analytical bounds on error probability for TCM on time-dispersive channels so that the issues of designing good codes on such channels can be studied.

Analytical upper bounds on error probability that rely on knowing the probability density function (pdf) of the ISI are derived. These bounds are closed-form expressions, but numerical techniques must be used to evaluate them. The emphasis in this work has been to obtain upper bounds that are tight for a wide range of time-dispersive channel conditions. A lower bound is also presented; this bound is tight for low levels of ISI, but loose for severe ISI.

The pdf of the ISI must be computed to evaluate the analytical upper bounds. However, exact computation of the pdf is only tractable in a few special cases. Algorithms are presented for computing approximations to the ISI pdf for uncoded and trellis-coded systems with general one- and two-dimensional signal constellations. Procedures for forming worst and best case ISI pdf's that can be used to compute upper and lower bounds on symbol error probability are also developed. Examples show that the dependence between symbols, introduced into the transmitted signal by Ungerboeck codes, has negligible effect on the pdf of the ISI.

In summary, the performance of TCM in a residual ISI environment has been studied using computer simulation. To overcome the limitations of simulation, analytical bounds on error probability have been developed. The numerical evaluation of these analytical bounds relies on algorithms to compute approximate ISI pdf's. The tightness of the bounds has been verified by computer simulation where possible.

Acknowledgements

There are a number of people whom I wish to thank for their contribution to my work. Mr. Bill Kennedy, from the University of Canterbury, has acted as one of my two co-supervisors and has carefully edited successive drafts of research papers and this thesis. Without Bill's efforts, my work would have been less clearly expressed. Dr. Mansoor Shafi, from Telecom Corporation of New Zealand Ltd., was my other co-supervisor and it was he who introduced me to the concept of trellis-coded modulation. He has instilled in me a deep enthusiasm for digital communication.

I am indebted to Professor Des Taylor, from McMaster University, for showing an interest in my work and for providing me with the opportunity and funding to spend five months working at McMaster University. Professor Taylor also assisted with the sponsorship that enabled me to attend ICC'89 and the 1989 IEEE Information Theory Workshop.

The financial support of Telecom Corporation of New Zealand Ltd. has been greatly appreciated. I am also grateful for the flexibility that allowed me to work in Canada and for sponsorship to attend ICC'89.

The postgraduate students and staff with whom I have been associated at the University of Canterbury have provided a friendly and stimulating environment in which to work. I would especially like to thank Martin Clark for the enlightening technical and philosophical discussions that we have had on digital communication.

Finally, I would like to thank the members of my family for their support; particularly my wife, Wendy, who proofread this thesis and tolerated my work habits over the last six months.

Contents

Preface	xi
Glossary	xv
Chapter 1 Introduction	1
1.1 Fundamentals of Digital Communication	2
1.1.1 Channel	3
1.1.2 Source Encoder and Decoder	5
1.1.3 Channel Encoder and Decoder	6
1.1.4 Design Goals and Constraints	7
1.2 Aim of the Thesis	8
Chapter 2 Some Mathematical Preliminaries	9
2.1 Geometric Representation of Signals	9
2.2 Representation of Bandpass Signals	10
2.3 Analysis of Bandpass Systems	11
2.3.1 Transmitter	11
2.3.2 Receiver	12
2.3.3 Performance of the Receiver	15
2.4 Conclusion	18
Chapter 3 Trellis-Coded Modulation	19
3.1 Approaching Shannon's Bound	20
3.2 Ungerboeck Encoding	22
3.2.1 Convolutional Encoding	22
3.2.2 Signal Mapping	23
3.3 Ungerboeck Decoding	25
3.3.1 Viterbi Decoding	26
3.4 Performance of Trellis-Coded Modulation	27
3.5 Conclusion	28
Chapter 4 Digital Microwave Radio Systems	29
4.1 Modulation and Demodulation	30
4.1.1 Modulation	30
4.1.2 Demodulation	31

4.1.3	Pulse Shaping	33
4.2	Multipath Fading	34
4.2.1	Channel Models	34
4.3	Countermeasures for Multipath Fading	37
4.3.1	Adaptive Equalization	37
4.3.2	Diversity	41
4.3.3	Error-Control Coding	41
4.4	Measures of Performance	42
4.5	Conclusion	43
Chapter 5	Performance Estimation by Simulation	45
5.1	Simulation Specifications	46
5.2	Error Rate Results	49
5.2.1	Additive White Gaussian Noise	50
5.2.2	Centred Spectral Notch	51
5.2.3	Offset Spectral Notch	52
5.2.4	Spectral Slope	54
5.2.5	Discussion of Results	55
5.3	Outage Probability Results	55
5.3.1	System Signatures	56
5.3.2	Outage Computation from Signatures	57
5.3.3	Error-Free Seconds and Residual Bit Error Rate	59
5.4	Conclusion	60
Appendix 5A	Cursor Attenuation	60
Chapter 6	The Probability Density of Intersymbol Interference	63
6.1	Preliminaries	64
6.2	Uncoded Systems	66
6.3	Trellis-Coded Systems	68
6.4	Best and Worse Case Binning	71
6.5	Examples of Probability Density Functions	72
6.5.1	Binning Parameters	73
6.5.2	Characteristics of the Probability Density Functions	75
6.6	Conclusion	78
Appendix 6A	An Upper Bound on Error Probability	79
Chapter 7	Analytical Performance Bounds	83
7.1	Preliminaries	84
7.2	Union Bound on Error-Event Probability	85
7.3	Upper Bounds on Conditional Pairwise Error Probability	86
7.3.1	Viterbi Bound	87
7.3.2	Chernoff Bound	89
7.4	Numerical Evaluation of the Union Bound	90
7.5	Numerical Evaluation of the Minimum Distance	93
7.6	Lower Bound on Error-Event Probability	93

7.7	Examples of Bounds on Error Probability	94
7.8	Conclusion	99
Appendix 7A	ISI Limit for Viterbi Bound	100
Chapter 8	Conclusions	101
8.1	Suggestions for Further Research	102
	References	105

Preface

Trellis-coded modulation (TCM) combines the functions of channel coding and modulation in a digital communication system. The redundancy introduced with TCM is accommodated in an expanded signal constellation so the system performance can be improved without sacrificing data rate or requiring increased bandwidth. I was introduced to the concept of TCM by Dr. Mansoor Shafi in September 1987. At this time, TCM had been predominantly considered for applications to additive white Gaussian noise (AWGN) channels.

The major impairment to transmission on digital microwave radio (DMR) systems is frequency selective fading, which can introduce severe intersymbol interference (ISI). Dr. Shafi proposed a research project to investigate whether or not TCM could combat the effects of the residual ISI that remains after non-ideal equalization on conventional uncoded DMR systems. We decided initially to perform a simulation study to see if the proposal was viable, and to follow this up with an analytical study to more generally express the performance of TCM. This thesis is the result of the research project and fulfills the aims of the original specification.

From May 1989 to November 1989 I worked with Professor Des Taylor and his communication research group at McMaster University, Ontario, Canada. Professor Taylor had visited Canterbury University under the Erskine fellowship program in July 1988 and our common research interests led to an invitation for me to visit McMaster University. During my stay in Canada I developed a significant amount of the analytical work contained in this thesis. I was also able to attend the 1989 IEEE International Conference on Communications (ICC'89) in Boston and the 1989 IEEE Workshop on Information Theory at Cornell University.

Chapter 1 briefly sketches the history of communications and, more specifically, describes some of the fundamental concepts of digital communication. The impact of Shannon's information theory on communication system design is highlighted. Finally, this chapter provides the motivation for the research in this thesis.

Some mathematical background to digital communications is provided in Chapter 2. A geometrical representation of signals is presented and used to illustrate the functions of the transmitter and receiver. The performance of the receiver, in terms of symbol error probability, is discussed.

Trellis-coded modulation and Digital Microwave Radio (DMR) systems are central to the research in this thesis. The principles of TCM are described in Chapter 3 and set in the broader context of channel coding. Relevant aspects of DMR systems are discussed in Chapter 4.

Chapter 5 contains details and results of a simulation study to examine the performance of DMR systems incorporating TCM. The simulation results verified that TCM will reduce the effects of residual ISI and gave motivation for an analytical investigation.

The next two chapters deal with analytical techniques. Chapter 6 examines the probability density of ISI. Algorithms to calculate approximate ISI probability density functions (pdf's) for uncoded and trellis-coded systems are derived. Modified forms of these algorithms are used to compute worst case and best case ISI pdf's. These pdf's can be used to obtain analytical lower and upper bounds on the probability of the system making a decision error.

Chapter 7 develops a union bound of pairwise error probabilities, using knowledge of the ISI pdf. The pairwise error probabilities must themselves be upper bounded so that a generalization of the transfer function of a convolutional code can be used to evaluate the bound. The union bounds are shown to be sufficiently tight to be useful in practice. A lower bound with ISI is given by the lower bound for TCM on an AWGN channel. This bound is tight for low levels of ISI, but is loose for severe ISI.

The thesis concludes with Chapter 8, which provides a summary of the results, and identifies outstanding issues that suggest ideas for further research.

The original research in this thesis is contained in Chapters 5, 6, and 7. A computer program was developed specifically for the simulation study. This study shows that TCM can offer significant improvements in performance for DMR systems with equalization. To my knowledge, no other simulation studies of DMR systems with TCM had been published at the time this study was undertaken.

Knowledge of the ISI pdf for a system on a time-dispersive channel is useful to estimate system error probabilities analytically. If an ISI pdf cannot be computed exactly, it is usually assumed to have the form of a standard pdf, such as a uniform pdf or a Gaussian pdf. The idea of approximating ISI pdf's for uncoded systems using quantization and binning procedures has been examined by Hill [1971] and Metzger [1987], but the algorithm presented in this thesis for approximating ISI pdf's, when the transmitted signal is trellis-coded, is original. The ideas of worst and best case binning for ISI pdf's are also original.

Few analytical techniques have previously been described for computing bounds on the error probability of TCM on time-dispersive channels. The techniques that have been described are restricted to analyzing simple systems and channels because they explicitly analyze the channel states. In this thesis, a union bound on error probability is developed specifically to make use of the availability of an approximation to the ISI pdf and avoid having to explicitly analyze the channel states. Thus, this union bound is applicable to a wide range of channels and systems. Generalized code transfer functions have been used previously to evaluate union bounds; however, the formulation of the union bound to account for ISI and multiplicative interference is largely original, although it is similar to Divsalar's analysis for mismatched receivers [Divsalar, 1978].

Six papers have been written as a result of this research, and are listed below. I presented the NELCON'88 paper in Christchurch in September 1988 and the ICC'89 paper in Boston, Massachusetts in June 1989. Ross McKay presented the GLOBECOM'88 paper in Hollywood, Florida in December 1988.

Carlisle, C.J., Shafi, M. and Kennedy, W.K. (1988), 'Trellis-coded modulation—approaching Shannon's bound', in *NELCON'88 Conf. Proc.*, Christchurch, N.Z., pp. 182–187.

Carlisle, C.J., Kennedy, W.K. and Shafi, M. (1989), 'Outage simulations for digital microwave radio systems with trellis-coded modulation', in *ICC'89 Conf. Rec.*, Boston, Mass., pp. 33.2.1–33.2.5.

Carlisle, C.J., Shafi, M. and Kennedy, W.K. (1990a), 'Trellis-coded modulation on digital

microwave radio systems—Simulations for multipath fading channels', accepted for publication in *IEEE Trans. Commun.*

Carlisle, C.J., Taylor, D.P., Kennedy, W.K. and Shafi, M. (1990b), 'The probability density of intersymbol interference for trellis-coded modulation', submitted to *IEEE Trans. Commun.*

Carlisle, C.J., Taylor, D.P., Kennedy, W.K. and Shafi, M. (1990c), 'Performance bounds for trellis-coded modulation on time-dispersive channels', submitted to *IEEE Trans. Commun.*

McKay, R.G., Shafi, M. and Carlisle, C.J. (1988), 'Trellis-coded modulation on digital microwave radio systems—Simulations for multipath fading channels', in *GLOBE-COM'88 Conf. Rec.*, Hollywood, Fla., pp. 8.1.1–8.1.5.

Glossary

Mathematical notation and abbreviations used in this thesis are defined below.

Mathematical Notation

The following definitions of mathematical notation are used. Most of the notation is standard, but it is included for completeness. Variables are not included here, but are defined when they are used in the thesis.

		Typical Reference
\triangleq	defined as	(2.6)
\equiv	equivalent to	(2.13)
\in	an element of (the set)	(1.3)
\iff	if and only if	(2.22)
:	such that	(6.25)
*	convolution	(6.6)
\sum'	summation with deletion of the zeroth term	(4.10)
j	$\sqrt{-1}$	(2.6)
$\text{Re}[z]$	real part of the complex number z	(2.7)
$\text{Im}[z]$	imaginary part of the complex number z	(4.10)
$ z $	magnitude of z	
$\angle z$	phase of z	(4.7)
z^*	complex conjugate of z	p. 33
A	matrix	(7.40)
A^t	matrix transpose	(7.40)
A^{-1}	matrix inverse	(7.42)
$\{A\}$	set of objects or events	p. 9
$P[A]$	probability of event A	(2.11)
$P[A B]$	probability of event A conditioned on event B	(2.11)
$p(x)$	probability density function of random variable X	(2.11)
$p(x y)$	probability density function of random variable X conditioned on Y	(2.16)
$E[X]$	expected value of the random variable X	(2.33)
\max	maximum element of a set	(7.56)
\min	minimum element of a set	(3.4)
\lim	in the limit	(7.56)
O	of the order of	p. 65

Abbreviations

The following abbreviations are also defined at the beginning of each chapter in which they are used.

		Typical Reference
AGC	automatic gain control	p. 37
AWGN	additive white Gaussian noise	p. 13
BER	bit error rate	p. 15
bits/s	bits per second	p. 6
bits/s/Hz	bits per second per Hertz	p. 20
BPSK	binary phase-shift keying	p. 2
CCITT	International Telephone and Telegraph Consultative Committee	p. 42
cm	centimetres	p. 1
CSS	coded signal set	p. 69
dB	decibels	p. 20
dB/MHz	decibels per MegaHertz	p. 50
DFE	decision-feedback equalizer	p. 39
DMR	digital microwave radio	p. 8
EFS	error-free seconds	p. 59
FSE	fractionally-spaced equalizer	p. 39
Gbits/s	Gigabits per second	p. 29
GHz	GigaHertz	p. 1
I	in-phase	p. 30
IF	intermediate frequency	p. 31
ISDN	integrated services digital network	p. 2
ISI	intersymbol interference	p. 4
kbits/s	kilobits per second	p. 2
kHz	kiloHertz	p. 1
km	kilometres	p. 1
m/s	metres per second	p. 1
MAP	maximum a posteriori probability	p. 13
Mbits/s	Mega bits per second	p. 29
MHz	MegaHertz	p. 46
ML	maximum-likelihood	p. 13
MLSE	maximum-likelihood sequence estimation	p. 15
MMSE	minimum mean-square error	p. 39
ns	nanoseconds	p. 37
OM	orders of magnitude	p. 50
pdf	probability density function	p. 12

PAM	pulse amplitude modulation	p. 11
PSK	phase-shift keying	p. 12
Q	quadrature	p. 30
QAM	quadrature amplitude modulation	p. 11
RF	radio frequency	p. 31
SER	symbol error rate	p. 15
SNR	signal-to-noise ratio	p. 4
SSE	synchronously-spaced equalizer	p. 38
TCM	trellis-coded modulation	p. 7
USS	uncoded signal set	p. 69
VLSI	very large-scale integration	p. 41
ZF	zero-forcing	p. 38

Chapter 1

Introduction

The ability of humans to communicate within shouting or visual range has not satisfied their need to exchange information. The more widely people travel, the greater the distances over which they wish to communicate. Some of the early systems for long distance communication were the use of smoke signals by the American Indians, the use of messengers, and signalling between a succession of towers using lanterns—much like modern communication systems use repeaters.

Samuel Morse is credited with developing the first widely successful electric telegraph on a wire circuit during the 1830s. The code he used is now known as the Morse code, which represents letters of the alphabet by spaces, dots, and dashes. A space was represented by an absence of electric current, a dot was a short duration pulse of current, and a dash was a longer duration pulse of current. Telegraphy has the disadvantage of not using a natural mode of communication and therefore not being directly accessible to most people. The first successful electric telephone was developed by Alexander Graham Bell in 1876. Although others had previously succeeded in using electricity to transmit sounds, Bell was the first to patent a device capable of sending and receiving recognizable words. The next step in the quest for long distance communication was to remove the restriction of wire circuits. Radio telegraphy was developed by Marconi, who succeeded in transmitting signals across the Atlantic in 1901. Speech was first transmitted by radio waves when, in 1906, Reginald Fessenden successfully applied the idea of modulating radio waves with a speech signal.

Modern communication systems are almost entirely based on the propagation of electromagnetic waves through a channel. Some common channels are the atmosphere, metallic conductors, optical fibres, and media for data storage. Electrical communication systems transmit information by using modulation to vary the amplitude, phase, or frequency of an electromagnetic wave. Systems that transmit electromagnetic waves in a band of frequencies around zero are known as *baseband systems*, and systems that transmit in a band of frequencies around a non-zero—carrier—frequency are known as *bandpass systems*.

In radio transmission, the carrier is converted to an electromagnetic field for propagation. To efficiently couple transmitted electromagnetic energy into space, the dimensions of the antenna aperture should be of the order of the wavelength $\lambda = c/f$ of the carrier, where f is the frequency of the carrier and $c = 3 \times 10^8$ m/s is the speed of light. If a baseband signal of $f = 3$ kHz is to be transmitted through space without modulating a carrier, the dimension of the antenna must be about 100 km for efficient coupling. In contrast, if the signal modulates a 30 GHz carrier, the required antenna aperture is only 1 cm. This illustrates the usefulness of modulation. Another reason for modulating onto

a carrier is that frequency division multiplexing can be used to allow many signals to be slotted into adjacent frequency bands.

A digital signal represents information with discrete symbols from a finite alphabet (i.e. discrete in time and amplitude). Therefore, the early telegraph systems like Morse's were digital systems. An analog signal is a function defined over a continuous range of time in which amplitude can assume a continuous range of values. The communication systems required to transmit digital signals are fundamentally different to those required to transmit analog signals.

The main classes of long distance communication traffic can be identified as person-to-person (speech or video), broadcast (radio or television), and inter-computer (data) communications. The first two of these classes use a predominantly analog format and the last one uses an exclusively digital format (analog computers aside). These different types of traffic are currently carried on separate communication networks. This scenario is rapidly changing, with the proposed Integrated Services Digital Network (ISDN) treating all communication traffic in a digital format (speech and video will be converted from analog format to digital format).

This thesis will deal only with digital communication systems.

1.1 Fundamentals of Digital Communication

It is ironic that, before the telephone was invented, all existing communication systems were digital, and now, after a century of analog domination, digital communication systems are again becoming dominant. There are a number of reasons for the use of digital communication systems; one of the most important reasons is that digital signals are much easier to regenerate than analog signals. Therefore, strategically placed repeaters allow the transmission of a digital signal over large distances, with significantly less distortion than an equivalent analog signal. Digital systems can use *time division multiple access* to allow many users access to a system. This is simpler than the *frequency division multiple access* that can be used in digital systems but must be used in analog systems. Digital systems can also use the spread spectrum technique of *code division multiple access* that promises to allow many more users to access a system simultaneously than either of the above access techniques [Schilling *et al.*, 1990]. The greater flexibility and reliability of digital hardware compared to analog hardware are further reasons for the current popularity of digital systems. Finally, digital systems can treat data, speech, and video signals identically, resulting in greater system flexibility and lower costs. This feature is used in the ISDN.

The major disadvantage of simple digital transmission systems is that, to communicate the same information, they generally require a greater system bandwidth than analog transmission systems. The bandwidth of an analog speech signal is about 4 kHz. For similar quality using *pulse code modulation* directly, the signal must be sampled at 8 kHz and each sample must be coded with 8 bits to give a 64 kbits/s digital signal. To transmit this signal using *binary phase-shift keying* (BPSK), a minimum bandwidth of 32 kHz is required. However, systems are now available that use source coding to reduce the data rate to 16 kbits/s, without significant loss in quality. But even this signal, transmitted using BPSK, requires a minimum bandwidth of 8 kHz—twice the analog bandwidth. Radio bandwidth is expensive; therefore, there is great motivation to use it efficiently. The efficient use of bandwidth can be achieved by the use of multilevel digital modulation.

The model of a general digital communication system is shown in Figure 1.1. The

transmitter consists of source and channel encoders, and processes a signal from the source, which may be analog or digital (e.g. voice or data), into a form suitable for the channel, which may also be analog or digital. The receiver consists of source and channel decoders, and processes the signal received on the channel to generate the best possible replica of the source, according to some fidelity criterion (usually minimum error probability).

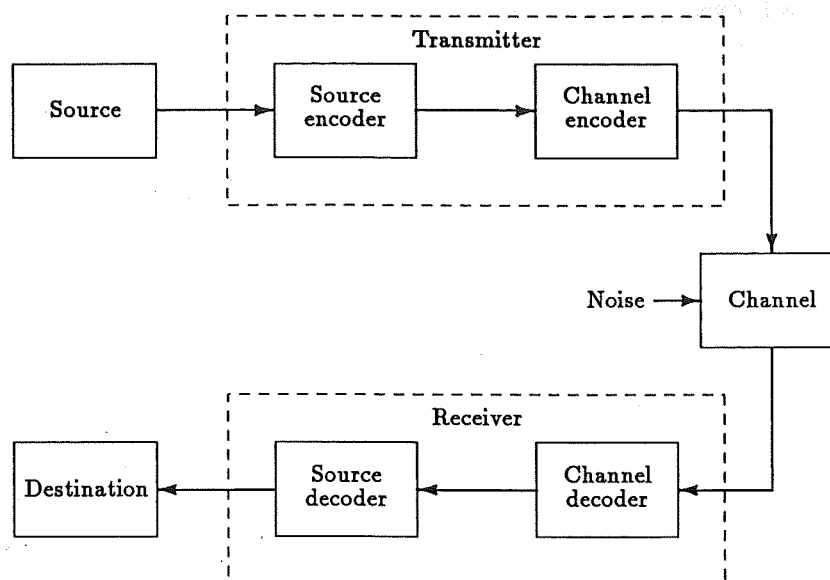


Figure 1.1 General digital communication system.

Two of the major contributors to the fundamentals of digital communication were Harry Nyquist and Claude Shannon. Nyquist [1928] showed that there is a fundamental lower limit on the bandwidth of the channel that can transmit pulses at a given signalling rate without mutual interference. Shannon laid the mathematical foundations for communication theory in 1948 with the publication of his paper: 'A Mathematical Theory of Communication' [Shannon, 1948a; 1948b]. This paper also created an entirely new field of endeavour, now known as *information theory*. Two of Shannon's most important theorems are known as the source coding theorem and the channel coding theorem. A very readable account of information theory is given by Pierce [1980]. The contributions of Nyquist and Shannon are discussed in this section.

1.1.1 Channel

If a signal could be transmitted over a channel and arrive undisturbed at the receiver, then the channel would be *ideal* and communication theory would be simple. However, such *ideal* channels do not exist in practice. There are a number of disturbances that occur on real channels and these introduce the possibility of errors when the receiver makes decisions on what was transmitted. Two of the most common disturbances are noise and intersymbol interference (ISI).

The ability to make correct decisions about the pulses transmitted is always limited by Gaussian noise, which is the result of thermal agitation of electrons in electronic circuits. Additive white Gaussian noise (AWGN) is encountered in all electrical communication systems. The effects of AWGN can be overcome by increasing the power of the transmitted pulses, but there are limits to the transmit power—governed by the cost of circuits, nonlinear effects, and government regulations.

AWGN is not always the dominant disturbance on a channel. Another disturbance, first experienced with early telegraph systems like Morse's (particularly with underground or undersea circuits, due to their high parallel capacitance), is that if a square pulse is transmitted it will be received as a spread out pulse. This phenomenon is illustrated in Figure 1.2. The speed of transmission is thus limited, because if pulses are transmitted too rapidly they will overlap in the receiver so that the receiver cannot sample one pulse without sampling part of another. Thus, the received pulses interfere with each other and decision errors result. This spreading of pulses so that they interfere with each other is known as ISI and is the predominant cause of error in many digital communication systems.

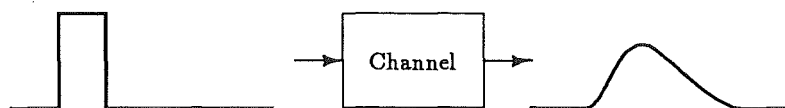


Figure 1.2 Spreading of a square pulse by the channel.

ISI occurs whenever a digital signal is passed through a linear channel with a frequency response that is not constant over the bandwidth of the signal. Furthermore, to produce a bandwidth efficient system, the bandwidth of the transmitted signal must be limited. The aim, therefore, is to design pulse shapes that occupy a minimum amount of bandwidth, and which result in a minimum of potential ISI.

Nyquist [1928] was one of the first researchers to recognize overlapping pulses as a source of interference. He is credited with formulating the criterion for zero ISI—the Nyquist criterion—which states that the ISI is zero if and only if T -spaced samples of the received pulse $h(t)$ satisfy

$$h(iT) = \begin{cases} 1 & \text{for } i = 0 \\ 0 & \text{for } i = \pm 1, \pm 2, \dots \end{cases} \quad (1.1)$$

for a signalling rate $1/T$ and where the sampling is synchronized with the pulse. In the frequency domain, this condition becomes

$$H_s(f) = \frac{1}{T} \sum_{j=-\infty}^{\infty} H(f - j/T) = K \quad \text{for } |f| \leq 1/2T \quad (1.2)$$

where K is a real constant, $H(f)$ is the channel frequency response, and $H_s(f)$ is the folded or aliased channel frequency response after symbol-rate sampling. The band of frequencies $|f| \leq 1/2T$ is known as the Nyquist or minimum bandwidth. An example of two successive Nyquist pulses is shown in Figure 1.3a. Notice that if the receiver samples at times iT , there will be no ISI.

Assuming that a pulse has been designed to satisfy the Nyquist criterion at a signalling rate $1/T$, ISI can be introduced by time-dispersion of the pulses in the channel or by modem imperfections (e.g. timing and carrier recovery errors, jitter, and imperfect filter design).

The narrower the pulse bandwidth, the larger are the tails or side lobes of the pulse, and the more severe the potential ISI. Bandwidth efficient systems are designed with pulse bandwidths close to the Nyquist bandwidth (some excess bandwidth is used to limit the severity of ISI due to channel distortion). ISI is often the main source of distortion for such systems operating at high *signal-to-noise ratios* (SNRs), and ultimately limits the system performance.

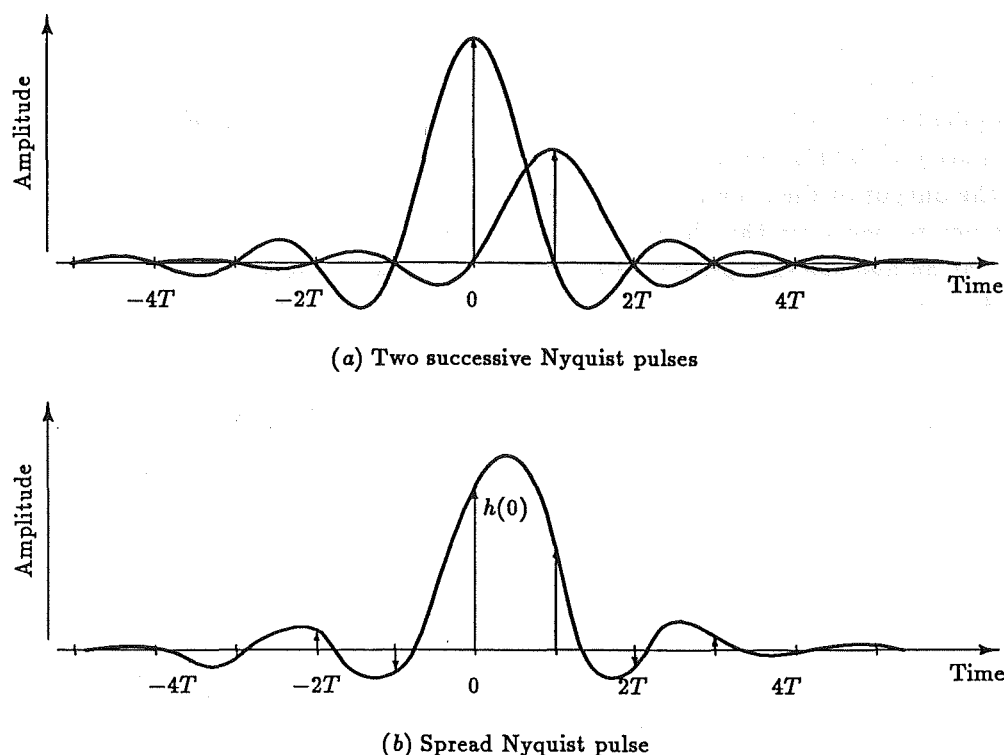


Figure 1.3 Examples of band-limited pulses.

A spread pulse is shown in Figure 1.3b. This pulse does not satisfy Nyquist's criterion for zero ISI. The main sample $h(0)$ is called the *cursor* of the sampled pulse, the samples $h(iT)$ for $i < 0$ are called *precursors*, and the samples $h(iT)$ for $i > 0$ are called *postcursors*.

ISI is typically mitigated in the receiver with an *equalizer* that attempts to cancel the effects of ISI by performing an inverse operation to the channel. If the channel is time varying (non-stationary), the equalizer can be designed to adapt to the changing channel conditions.

Nyquist [1928] also showed that the maximum number of distinct symbols that can be sent over a channel per second is twice the total bandwidth of frequencies used. Thus, the rate that symbols can be transmitted—the signalling rate—is proportional to bandwidth. The dual to this theorem is the Nyquist sampling theorem, which states that an analog signal with highest frequency component W can be represented completely and reconstructed perfectly over a time T seconds from a set of $2WT$ samples of its amplitude spaced $1/2W$ seconds apart.

1.1.2 Source Encoder and Decoder

Following the work of Nyquist, Shannon [1948a] defined a logarithmic measure of information. At the suggestion of J. W. Tukey, he called the unit of information a 'bit'—a contraction of 'binary digit'. Shannon also defined the entropy of a random source U , with specific messages $u \in \mathcal{U}$, as

$$H(U) = - \sum_{u \in \mathcal{U}} P[U = u] \log_2 P[U = u] \quad (1.3)$$

which is a measure of the average information, in bits per symbol, contained in messages produced by the source. If the source is analog, we assume it is digitized prior to further processing.

Suppose that the channel, channel encoder, and channel decoder form an error-free (noiseless) channel of finite capacity. Shannon's fundamental theorem for the noiseless channel [Shannon, 1948a] can be expressed as follows. Consider a noiseless channel with finite capacity C bits/s and a source with entropy H bits per symbol. It is possible to encode the output of the source in such a way as to transmit at the average rate $C/H - \epsilon$ symbols per second over the channel, where ϵ is arbitrarily small. It is not possible to transmit at an average rate greater than C/H . This theorem is also known as the *source coding theorem*.

The source coding theorem suggests the removal of redundant information from the source, using source coding, so that the encoded source rate R does not exceed the capacity of the channel ($R \leq C$). The encoded source rate cannot be less than the rate of entropy of the source H/T , without losing information. Morse code is an early example of source coding, where the most common letters of the alphabet are represented by the shortest codewords. If the source encoder mapping is one-to-one, the source decoder can simply perform the inverse mapping and deliver an exact replica of the source to the destination. When the source is analog, however, it cannot be represented perfectly by a digital sequence because the amplitude samples of the analog signal take on a continuum of values. In this case some distortion must be tolerated at the destination because the source decoder can only approximate the inverse mapping.

1.1.3 Channel Encoder and Decoder

Shannon [1948a] formulated an upper bound for the rate that information could be sent, with arbitrarily low error probability, over a band-limited additive white Gaussian noise (AWGN) channel. He used a geometrical representation of the channel to prove that the bound is exact [Shannon, 1949]. The upper bound is commonly known as the channel capacity

$$C = W \log_2(1 + S/N) \quad (1.4)$$

where W is the channel bandwidth, S is the signal power, and N is the noise power. This can be rearranged to give a theoretical limit on bandwidth efficiency

$$\zeta = C/W = \log_2(1 + S/N) \quad (1.5)$$

Shannon's fundamental result for the noisy channel [Shannon, 1948a] can be expressed as follows. If the rate of entropy of the source does not exceed the capacity of a Gaussian noise channel ($H/T \leq C$), then messages from the source can be transmitted over the channel with an arbitrarily small probability of error $P[\mathcal{E}]$. It is not possible to achieve an arbitrarily small probability of error if $H/T > C$. This theorem is also known as the *channel coding theorem*.

The channel coding theorem implies the addition of redundancy, using a channel encoder, to achieve arbitrarily small error probability, and as such is a dual to the source coding theorem. The goal of the channel encoder and decoder is to map the input data symbols into channel input symbols and conversely the channel output symbols into output data symbols, such that the effect of the channel noise is minimized.

Shannon showed that an arbitrarily small $P[\mathcal{E}]$ cannot be achieved at the channel capacity with any finite encoding process (finite delay and complexity), but it can be approached as closely as desired by using increasingly sophisticated coding schemes. His proof is an existence rather than a constructive proof, and he did not propose specific systems that could achieve an arbitrarily small $P[\mathcal{E}]$ at the channel capacity.

As a result of Shannon's theories we have seen the development of sophisticated channel coding techniques—as well as source coding techniques—in an effort to design bandwidth efficient systems and to approach the channel capacity.

Channel coding requires the transmission of redundant bits. These additional bits can be sent, without reducing the information rate, in one of two ways: either the signalling rate can be increased, if the channel bandwidth can be expanded, or the number of channel signals can be increased, if the channel is band-limited. However, the latter technique gives disappointing results when the coding and modulation operations are designed independently in the conventional manner.

Trellis-coded modulation (TCM) is a relatively new technique, developed in the 1970s by Ungerboeck [1982], that combines the functions of channel coding and modulation. The performance of a system can be improved by TCM without sacrificing data rate or requiring bandwidth expansion.

1.1.4 Design Goals and Constraints

The fundamentals of digital communication can be summarized by a list of goals and constraints that characterize the design of digital communication systems. There are a number of goals that we seek to achieve in designing a communication system:

1. Maximize the information rate R .
2. Minimize the probability of error $P[\mathcal{E}]$.
3. Minimize the transmit power S .
4. Minimize the required system bandwidth W .
5. Maximize the system use.
6. Minimize the system complexity, computational load, and system cost.

Most of these goals are in conflict with each other. Shannon, however, has shown us that 1 and 2 can be achieved independently, provided $R \leq C$. In seeking to achieve all the above goals, the designer faces a number of theoretical and regulatory constraints:

1. Nyquist theoretical minimum bandwidth requirement.
2. Source coding theorem.
3. Channel coding theorem.
4. Regulations (e.g. spectrum allocations).
5. Technological limitations.
6. Other system requirements (e.g. satellite orbits).

To cope with these constraints, one system parameter must be traded against another. For example, channel coding can be used to trade increased complexity (of the coding scheme) for one of the other goals. TCM can be used to achieve these trade-offs more efficiently than conventional coding techniques.

1.2 Aim of the Thesis

Most long distance transmission systems use either copper circuits, optical fibre, satellite repeaters, or terrestrial microwave radio. In this thesis we will concentrate on *digital microwave radio* (DMR) systems, although the techniques and results presented are applicable to other channels. DMR uses line-of-sight propagation and repeaters to transmit over long distances. The major impairment on DMR channels is sporadic multipath propagation. This results in multiple versions of the same signal, with different attenuation and delay, arriving at the receiver and interfering with each other. Multipath propagation is caused by inhomogeneous temperature and humidity profiles in the atmosphere. Fortunately, the rate of fading is slow compared to the signalling rate, and adaptive equalization can be used to reduce the severe ISI that results from deep fades.

We consider the performance of TCM on DMR systems. TCM is particularly suited to DMR systems because it can provide large coding gains without incurring a loss of bandwidth efficiency. Specifically, we study how TCM performs with *residual ISI* (i.e. the ISI that remains after non-ideal equalization), since TCM cannot be expected to cope with the *raw ISI* generated by a severe fade. Note that residual ISI could also include ISI due to modem imperfections, although we do not consider this situation. We wish to show that the performance of a trellis-coded system with residual ISI is superior to that of an equivalent uncoded system. This corresponds to trading code complexity for improved performance. TCM also allows code complexity to be traded for reduced transmit power or increased bandwidth efficiency.

The performance of TCM on an AWGN channel is well established in terms of lower and upper bounds on error probability and an asymptotic coding gain at high SNRs (see Section 3.4). These results extend directly to the DMR channel under normal (unfaded) propagation conditions and under flat fading conditions. The performance of TCM with time-dispersive fading is not so well defined, and we concentrate on this problem.

In this investigation, both simulation and analytical techniques have been used. The simulation study used Monte Carlo techniques whereby the system was simulated in computer software, data were transmitted over the system, and the frequency of errors in the received data was monitored. Computer simulations, however, are notoriously slow and limited in practice to high bit error rates (BERs). This inevitably leads to a search for analytical solutions. The most desirable analytical solution is a tractable closed-form expression for the error probability. Unfortunately, such expressions rarely exist, and analytical bounds on the error probability must be used. The approach taken here is to compute a good approximation to the probability density of the ISI in the receiver and to use this in the formulation of a union bound on the error event probability in the Viterbi decoder.

The simulation study is presented in Chapter 5. Due to the limitations of simulation, a major aim of this thesis has been to develop suitable analytical techniques. Algorithms for the computation of probability density functions of ISI are developed in Chapter 6, and examples are studied. The analytical bounds are formulated in Chapter 7. These bounds are computed for a number of examples and compared to the simulation results.

Chapter 2

Some Mathematical Preliminaries

The mathematical foundations of digital communication are well established. Wozencraft and Jacobs [1965] popularized the geometric representation of signals, which is discussed in Section 2.1, although the concept was first used by Shannon [1949]. Section 2.2 shows how bandpass signals and linear bandpass systems can be represented by equivalent low-pass forms to simplify the analysis of bandpass systems. Equivalent lowpass signals and geometric representations are used in Section 2.3 to discuss the operation of components of bandpass systems. This chapter does not provide a complete mathematical basis for communication theory, but merely introduces some basic terminology used in subsequent chapters. There are many books that give comprehensive accounts of communication theory, for example Wozencraft and Jacobs [1965], Viterbi and Omura [1979], and Proakis [1989].

2.1 Geometric Representation of Signals

An L -dimensional *orthogonal space* is defined by a set of L linearly independent *basis functions* $\{\varphi_j(t)\}$. Any arbitrary function in the space can be generated by a linear combination of these basis functions, provided they satisfy the orthogonality condition

$$\int_{-\infty}^{\infty} \varphi_j(t) \varphi_k(t) dt = \begin{cases} K_j & \text{if } j = k \\ 0 & \text{otherwise} \end{cases} \quad \text{for } j, k = 1, \dots, L \quad (2.1)$$

where the constants $\{K_j\}$ are nonzero. When the basis functions are normalized so each $K_j = 1$, the space is called *orthonormal*.

Any arbitrary set of M waveforms $\{y_i(t)\}$ can be expressed as a linear combination of L ($\leq M$) orthogonal waveforms $\{\varphi_j(t)\}$, such that

$$y_i(t) = \sum_{j=1}^L y_{ij} \varphi_j(t) \quad \text{for } i = 1, \dots, M \quad (2.2)$$

where

$$y_{ij} = \frac{1}{K_j} \int_{-\infty}^{\infty} y_i(t) \varphi_j(t) dt \quad \text{for } i = 1, \dots, M \text{ and } j = 1, \dots, L \quad (2.3)$$

The form of the basis functions $\{\varphi_j(t)\}$ is chosen for convenience and depends on the form of the signal waveforms. The waveform $y_i(t)$ can be viewed as a vector or *signal point* $\mathbf{y}_i = (y_{i1}, y_{i2}, \dots, y_{iL})$ in an L -dimensional *Euclidean space*. This is an extremely useful representation because, as we will show in Section 2.3.2, *Euclidean distance* is an optimum metric for the receiver if the channel only adds Gaussian noise to the transmitted signal.

Euclidean distance is a generalization of the concept of distance in three-dimensional physical space.

The set of M signal points in the signal space is called a *signal set* or a *signal constellation*. Baseband signals are inherently one-dimensional and can be represented by a scalar variable. Bandpass signals are inherently two-dimensional and can be represented by a complex variable, as discussed in the next section. Signals of higher dimensionality can be constructed using time- or frequency-orthogonal signals.

2.2 Representation of Bandpass Signals

A bandpass signal with carrier frequency f_c can be represented as

$$y(t) = a(t) \cos(2\pi f_c t + \theta(t)) \quad (2.4)$$

where $a(t)$ is the amplitude envelope of the signal and $\theta(t)$ is the phase of the signal. The cosine can be expanded to yield

$$y(t) = x_c(t) \cos(2\pi f_c t) - x_s(t) \sin(2\pi f_c t) \quad (2.5)$$

where the *in-phase* component $x_c(t) \triangleq a(t) \cos \theta(t)$ and the *quadrature* component $x_s(t) \triangleq a(t) \sin \theta(t)$ are lowpass signals. From these lowpass signals we can define a *complex envelope*

$$\begin{aligned} x(t) &\triangleq a(t)e^{j\theta(t)} \\ &= x_c(t) + jx_s(t) \end{aligned} \quad (2.6)$$

so that the bandpass signal can be represented as

$$y(t) = \operatorname{Re} [x(t)e^{j2\pi f_c t}] \quad (2.7)$$

The complex envelope defines all the characteristics of the bandpass signal, except the carrier. Since the carrier does not convey information, we can represent the bandpass signal by the *equivalent lowpass signal* $x(t)$. Similarly, a linear bandpass system can be represented by an equivalent lowpass impulse response $h(t)$ [Proakis, 1989].

A linear bandpass system can be analyzed with complex envelopes. When a bandpass signal with complex envelope $x(t)$ is applied to a bandpass system with complex impulse response $h(t)$, the complex envelope of the output can be expressed as [Proakis, 1989]

$$r(t) = \int_{-\infty}^{\infty} x(\tau)h(t - \tau) d\tau \quad (2.8)$$

Complex envelopes are universally used in the analysis of bandpass systems, and are particularly useful for simulation because the carrier need not be simulated. Simulating the carrier would require a sampling rate in excess of twice the carrier frequency. A *discrete Fourier transform* performed using a *fast Fourier transform* algorithm [Oppenheim and Schaffer, 1975] can be used to evaluate (2.8) numerically, leading to further efficiencies in computation.

2.3 Analysis of Bandpass Systems

The geometric representation of signals facilitates the use of equivalent discrete-time vector systems to analyze digital communication systems. Figure 2.1 shows a discrete-time communication system where the signals are two-dimensional vectors, represented by complex variables. The system transmits a random message U , with specific value $u \in \mathcal{U}$, by mapping it to a signal $X = f(U)$, with specific value $x \in \mathcal{X}$. The channel disturbs the signal to form the received signal R with specific value r from which a decision $\hat{u} \in \mathcal{U}$ must be made. Wozencraft and Jacobs [1965] discuss general L -dimensional systems.

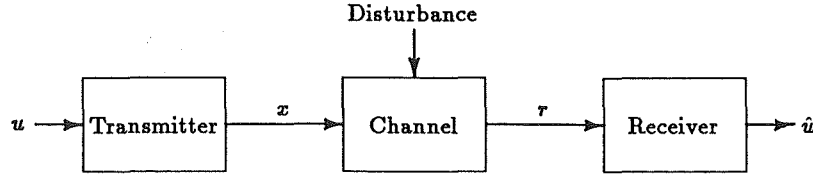


Figure 2.1 Discrete-time communication system.

The transmitter, receiver, and performance of the receiver are discussed in this section.

2.3.1 Transmitter

Consider a bandpass system that transmits a single pulse $x(t)$ with shape $h_T(t)$. The transmitted signal can be represented as

$$y_i(t) = \text{Re} \left[(A_{ic} + jA_{is}) h_T(t) e^{j2\pi f_c t} \right] \quad \text{for } i = 1, 2, \dots, M \quad (2.9)$$

where $\{x = A_{ic} + jA_{is}\}$ are M possible complex pulse amplitudes. If the pulse shape $h_T(t)$ is square and its duration is a multiple of the carrier period $1/f_c$ (or much greater than it), then basis functions for the signal space are $\varphi_1(t) = h_T(t) \cos(2\pi f_c t)$ and $\varphi_2(t) = h_T(t) \sin(2\pi f_c t)$, and $\{(A_{ic}, A_{is})\}$ are points in a two-dimensional signal space. The bandpass transmitter can thus be represented as in Figure 2.2, where the signal mapper is the equivalent lowpass transmitter.

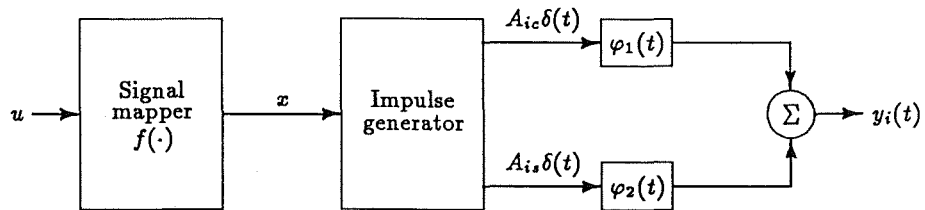


Figure 2.2 Bandpass transmitter.

When $A_{is} = 0$ for $i = 1, 2, \dots, M$, the signal points only occupy one dimension of the signal space, and the bandpass signal is called M -ary *pulse amplitude modulation* (M-PAM). A 4-PAM signal constellation is illustrated in Figure 2.3a.

Bandpass PAM signals use bandwidth inefficiently unless *single sideband* transmission is used [Proakis, 1989, Chapter 3]. This involves removing the redundant sideband, and requires expensive hardware. A more efficient bandpass signal is obtained when the M signal points $\{(A_{ic}, A_{is})\}$ form a rectangular grid, as shown in Figure 2.3b for $M = 16$. This modulation is called M -ary *quadrature amplitude modulation* (M-QAM).

When a system is severely nonlinear (e.g. when the *high power amplifier* is in saturation), the signal amplitude and the information contained within it are distorted.

The problem of amplitude distortion can be avoided by using M -ary *phase-shift keying* (M-PSK), which transmits information using only the phase of the carrier. This can be achieved by using signal points with constant amplitude $\sqrt{A_{ic}^2 + A_{is}^2}$, and phase $\tan^{-1}(A_{is}/A_{ic}) = 2\pi i/M$. An 8-PSK signal constellation is shown in Figure 2.3c.

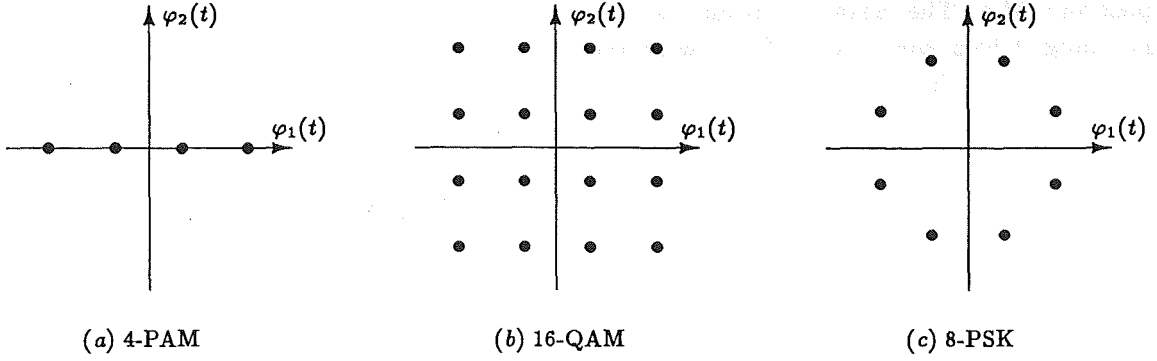


Figure 2.3 Examples of signal constellations.

If a sequence of symbols (pulses) is transmitted, the equivalent lowpass signal transmitted over the channel is

$$x(t) = \sum_{n=-\infty}^{\infty} x_n h_T(t - nT) \quad (2.10)$$

where $\{x_n\}$ is the sequence of complex symbols and $h_T(t)$ should satisfy the Nyquist criterion for zero ISI when sampled at $t = nT$ for integer n . The condition for zero ISI is essentially a time-orthogonality condition. Time-orthogonal pulses can be used to form signals with more than two dimensions. In an analogous manner, frequency-orthogonal signals can be used to add further dimensions.

2.3.2 Receiver

The most appropriate—but generally intractable—criterion for designing an optimum receiver is to minimize the probability of an erroneous symbol decision $P[\mathcal{E}]$, or equivalently to maximize the probability of a correct decision

$$P[C] = \int_{-\infty}^{\infty} P[C | R = r] p_R(r) dr \quad (2.11)$$

where $P[C | R = r]$ is the probability of a correct decision conditioned on R , and $p_R(r)$ is the *probability density function* (pdf) of R . From (2.11) we see that $P[C]$ is maximized when $P[C | R = r]$ is maximized. When the receiver sets $\hat{u} = u'$, the conditional probability of a correct decision is

$$P[C | R = r] = P[U = u' | R = r] \quad (2.12)$$

where $P[U = u' | R = r]$ is the *a posteriori probability* of message u' having been transmitted.

To avoid cumbersome notation in subsequent chapters, we use the equivalent notation

$$P[a] \equiv P[A = a] \quad \text{for all } a \in \mathcal{A} \quad (2.13)$$

where A is usually a discrete random variable with specific value $a \in \mathcal{A}$. We also use

$$p(b) \equiv p_B(b) \quad \text{for all } b \quad (2.14)$$

where B is usually a continuous random variable with specific value b . This notation will be used when there is no possibility of ambiguity.

The optimum receiver must determine, for a given received signal, which of the messages $u \in \mathcal{U}$ has *maximum a posteriori probability* (MAP). The so called MAP receiver uses its knowledge of the conditional pdf $p(r | x)$, the signal set \mathcal{X} , and the *a priori* message probabilities $P[u]$ to set $\hat{u} = u'$ whenever

$$P[u' | r] > P[u | r] \quad \text{for all } u \neq u' \text{ and } u, u' \in \mathcal{U} \quad (2.15)$$

If two or more messages have MAP, the receiver arbitrarily selects one of them. Using the mixed form of Bayes' rule [Wozencraft and Jacobs, 1965, Chapter 2]

$$P[u' | r] = \frac{P[u'] p(r | u')}{p(r)} \quad (2.16)$$

and

$$p(r | u') = p(r | x') \quad \text{for } u' \in \mathcal{U} \text{ and } x' = f(u') \in \mathcal{X} \quad (2.17)$$

where $f(\cdot)$ is the signal mapping function. Since $p(r)$ is independent of the transmitted message, the MAP receiver sets $\hat{u} = u'$ whenever the decision function

$$P[u'] p(r | x') \quad \text{for } u' \in \mathcal{U} \text{ and } x' = f(u') \in \mathcal{X} \quad (2.18)$$

is maximum.

The *a priori* message probabilities $P[u]$ are often unknown. In such cases, the receiver can determine \hat{u} to maximize $p(r | x)$. This receiver is called a *maximum-likelihood (ML) receiver*, and is optimum when all messages are equally likely.

If the disturbance on the channel is *additive white Gaussian noise* (AWGN), the received signal is $r = x + \eta$, where η is a specific value of a complex Gaussian random variable \mathcal{N} . The AWGN actually has an infinite number of dimensions, but only the vectors that lie within the signal space are relevant to the decision process [Wozencraft and Jacobs, 1965, Chapter 4]. The ML receiver determines \hat{u} by maximizing

$$p(r | x') = p_{\mathcal{N}}(r - x' | x') \quad (2.19)$$

The random variables \mathcal{N} and X are usually statistically independent, so

$$p_{\mathcal{N}}(r - x' | x') = p_{\mathcal{N}}(r - x') \quad \text{for all } x' \in \mathcal{X} \quad (2.20)$$

The pdf of the Gaussian noise is

$$p_{\mathcal{N}}(\alpha) = \frac{1}{2\pi\sigma_{\eta}^2} \exp\left(-\frac{|\alpha|^2}{2\sigma_{\eta}^2}\right) \quad (2.21)$$

with zero mean and variance σ_{η}^2 . Therefore, maximizing $p(r | x') = p_{\mathcal{N}}(r - x')$ is equivalent to minimizing $|r - x'|^2$. Geometrically, $|r - x'|^2$ is the squared Euclidean distance between points r and x' in the signal space. Thus the ML decision rule assigns the received point r to u' if and only if r is closer to $x' = f(u')$ than any other signal point. This decision rule can be used to partition the signal space into M disjoint regions $\{\Lambda_{u'}\}$ so that whenever the received vector r is in $\Lambda_{u'}$, the receiver sets $\hat{u} = u'$

$$r \in \Lambda_{u'} \iff \hat{u} = u' \quad (2.22)$$

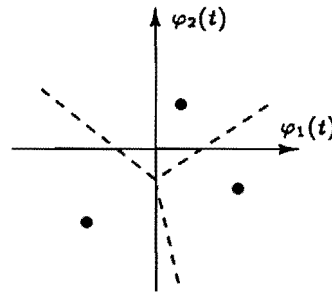
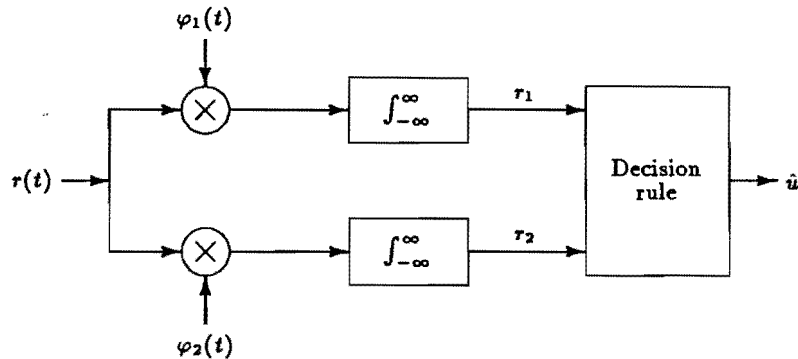


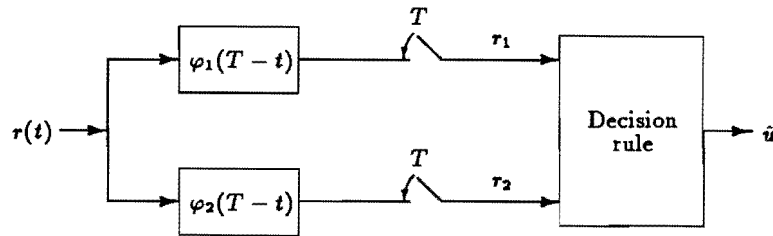
Figure 2.4 ML decision regions for three signal points.

These regions are called decision regions and are illustrated in Figure 2.4 for three equally likely signal points. For the ML decision rule, the decision boundary between any two signal points is the perpendicular bisector of a line between the points.

Before the receiver can apply a decision rule, it must extract a received signal vector from the received signal. In the absence of noise, the receiver can extract the received signal vector from $r(t)$ by correlating $r(t)$ with each of the basis functions $\varphi_j(t)$, as implied by (2.3); this forms the basis of the *correlation receiver* shown in Figure 2.5a. When noise is present, the correlation receiver maximizes the *signal-to-noise ratio* [Wozencraft and Jacobs, 1965, Chapter 4]. The correlators in the correlation receiver can equivalently be implemented by *matched filters* $\{\varphi_j(T - t)\}$ sampled at time T , to obtain the *matched filter receiver* shown in Figure 2.5b. The basis functions $\{\varphi_j(t)\}$ should be zero for $t > T$



(a) Correlation receiver



(b) Matched filter receiver

Figure 2.5 ML bandpass receivers.

so that the matched filters are physically realizable. After matched filtering, the relevant noise components in the received signal vector are independent Gaussian random variables. Generally demodulation is regarded as the process of removing the carrier from the received signal (obtaining the lowpass signal components) and sampled matched filtering is regarded as obtaining the time-orthogonal signal components, although both operations are actually

correlation operations.

A system that transmits a square pulse $h_T(t)$ requires an infinite bandwidth channel to avoid distortion, but in practice we are usually interested in band-limited or near band-limited (most of the pulse energy within a finite bandwidth) systems. Also, a system that transmits a single message is not very useful; therefore, we are usually interested in systems that can transmit sequences of messages. Both these issues can be addressed by designing band-limited pulses that satisfy Nyquist's criterion for zero ISI (time-orthogonality). Sequences of messages can be transmitted and received by time-sharing the transmitter and receiver over successive messages. If the Nyquist pulses are distorted on the channel, the pulses are no longer time-orthogonal and ISI results in the receiver.

When a channel introduces ISI in addition to AWGN, the linear receivers just described are no longer optimum. The optimum linear receiver becomes a sampled matched filter followed by an inverse filter (infinite-tap transversal equalizer) prior to applying the ML decision rule [Proakis, 1989, Chapter 6]. This receiver will enhance the Gaussian noise and is not globally optimum. An optimum nonlinear receiver that does not enhance the Gaussian noise is described by Forney [1972]. This receiver consists of a *whitened matched filter* followed by *maximum-likelihood sequence estimation*. The whitened matched filter consists of a matched filter followed by a transversal filter to whiten the noise. The noise must be whitened so that a squared Euclidean distance metric can be used with the Viterbi algorithm, to determine the maximum-likelihood transmitted sequence. The difficulty of implementing Forney's receiver has prompted a number of researchers to examine suboptimum, but realizable, receivers for ISI channels [Magee and Proakis, 1973; Qureshi and Newhall, 1973; Falconer and Magee, 1973; Ungerboeck, 1974; Eyuboğlu and Qureshi, 1988].

2.3.3 Performance of the Receiver

The decision error probability $P[\mathcal{E}]$ forms the basis of most performance measures. In practice, *bit error rate* (BER) or *symbol error rate* (SER) are used; these measures are defined as the ratio of the number of errored data units to total data units over a specified time interval.

The probability of a symbol decision error can be expressed as

$$P_s[\mathcal{E}] = \sum_{u \in \mathcal{U}} P[u] P_s[\mathcal{E} | u] \quad (2.23)$$

Defining $\{\bar{\Lambda}_u\}$ as the set of erroneous decision regions for messages $u \in \mathcal{U}$, the conditional probability of a symbol decision error is

$$P_s[\mathcal{E} | u] = P[r \in \bar{\Lambda}_u | u] \quad (2.24)$$

$$= \int_{\bar{\Lambda}_u} p(r | f(u)) dr \quad (2.25)$$

For an AWGN channel

$$P_s[\mathcal{E} | u] = \int_{\bar{\Lambda}_u} p_{\mathcal{N}}(r - f(u)) dr \quad (2.26)$$

$$= \frac{1}{2\pi\sigma_n^2} \int_{\bar{\Lambda}_u} e^{-|r-f(u)|^2/2\sigma_n^2} dr \quad (2.27)$$

For example, consider the performance of an uncoded system using an M -PAM constellation, with minimum distance d between the signal points, on an AWGN channel.

The decision regions for an ML receiver are illustrated in Figure 2.6 for $M = 8$. The two points $\pm(M-1)d/2$ have a conditional symbol error probability of

$$P_s[\mathcal{E} | u] = Q\left(\frac{d}{2\sigma_\eta^2}\right) \quad (2.28)$$

where

$$Q(s) = \int_s^\infty e^{-t^2/2} dt \quad (2.29)$$

is the Gaussian integral function, also referred to as the Q -function. The other $M-2$ signal points have

$$P_s[\mathcal{E} | u] = 2Q\left(\frac{d}{2\sigma_\eta^2}\right) \quad (2.30)$$

Using (2.23) with $P[u] = 1/M$ we get an expression for the probability of a symbol error

$$P_s[\mathcal{E}] = 2\frac{M-1}{M} Q\left(\frac{d}{2\sigma_\eta^2}\right) \quad (2.31)$$

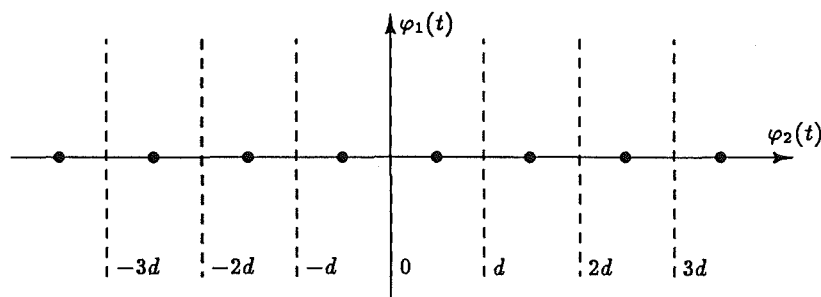


Figure 2.6 ML decision regions for M-PAM.

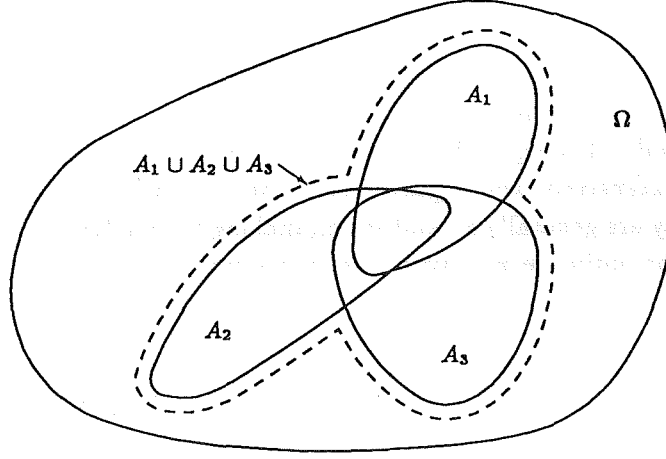
It is often not possible to form a tractable analytical expression for the decision error probability of a system, especially when the system uses channel coding or the channel introduces ISI. In these cases it is necessary to use probability bounds to estimate the probability of an error. In Chapter 7 we will use the Union, Chernoff, and Viterbi bounds that are introduced below.

Union Bound

Expressions for $P[\mathcal{E}]$ usually involve the probability of a union of events $P\left[\bigcup_{i=1}^K A_i\right]$, where $\{A_i\}$ is a set of events from the set Ω . This union is difficult to evaluate exactly, especially for large K . An examination of the Venn diagram in Figure 2.7 shows that, for $K = 3$,

$$P\left[\bigcup_{i=1}^K A_i\right] \leq \sum_{i=1}^K P[A_i] \quad (2.32)$$

This is true for any K , and the bound is known as the *union bound* [Wozencraft and Jacobs, 1965].

Figure 2.7 Venn diagram for $K = 3$ intersecting events.

Chernoff Bound

The *Chernoff bound* for a single random variable Y is derived by Proakis [1989] as

$$P[Y \geq \epsilon] \leq E[e^{\lambda(Y-\epsilon)}] \quad \lambda, \epsilon \geq 0 \quad (2.33)$$

where $E[\cdot]$ is the expected value and the Chernoff parameter λ is determined from

$$\frac{\partial}{\partial \lambda} E[e^{\lambda(Y-\epsilon)}] = 0 \quad (2.34)$$

Viterbi Bound

A bound that generally yields tighter bounds than the Chernoff bound is given by

$$Q(\sqrt{u+v}) \leq e^{-v/2} Q(\sqrt{u}) \quad u, v \geq 0 \quad (2.35)$$

We call this bound the *Viterbi bound* because Viterbi [1971] appears to be the first researcher to have applied it to digital communication problems. It can be proven by writing the Q -functions in terms of integrals. The right-hand side of (2.35) gives

$$\begin{aligned} e^{-v/2} Q(\sqrt{u}) &= e^{-v/2} \frac{1}{\sqrt{2\pi}} \int_{\sqrt{u}}^{\infty} e^{-s^2/2} ds \\ &= \frac{1}{\sqrt{2\pi}} \int_{\sqrt{u}}^{\infty} e^{-(v+s^2)/2} ds \end{aligned} \quad (2.36)$$

Using the substitution $v + s^2 = w^2$ and hence $ds = w dw / \sqrt{w^2 - v}$, we get

$$e^{-v/2} Q(\sqrt{u}) = \frac{1}{\sqrt{2\pi}} \int_{\sqrt{u+v}}^{\infty} \frac{w}{\sqrt{w^2 - v}} e^{-w^2/2} dw \quad (2.37)$$

Since $v \geq 0$, it follows that $w \geq 0$ and $w/\sqrt{w^2 - v} \geq 1$. Using these inequalities, we can write the inequality

$$\begin{aligned} \frac{1}{\sqrt{2\pi}} \int_{\sqrt{u+v}}^{\infty} e^{-t^2/2} dt &\leq \frac{1}{\sqrt{2\pi}} \int_{\sqrt{u+v}}^{\infty} \frac{w}{\sqrt{w^2 - v}} e^{-w^2/2} dw \\ &\leq e^{-v/2} Q(\sqrt{u}) \end{aligned} \quad (2.38)$$

thus proving the bound in (2.35).

2.4 Conclusion

The mathematical background presented in this chapter establishes some important results on which the rest of this thesis can build. Extensive use is made of the geometrical representation of signals. The equivalent lowpass representation for bandpass signals and systems is also used extensively, and is particularly useful for the simulations in Chapter 5. The receivers we study are generally suboptimum, making the analysis of their performance more difficult than for optimum receivers. Hence we must use bounds to study receiver performance analytically.

Chapter 3

Trellis-Coded Modulation

In a conventional data transmission system, the channel encoder and signal mapping function of the modulator are designed independently, so that the channel coding is essentially an addition to an uncoded system. The code redundancy is accommodated by reducing the data rate or by increasing the system bandwidth. In either case, the bandwidth efficiency is reduced relative to the uncoded system. The modulator maps m -bit code symbols into one of 2^m possible channel symbols. *Gray encoding* is usually used so that neighbouring signal points in the constellation are assigned to code symbols that differ in only one bit. Hence, the most likely symbol decision errors only result in single bit errors, and the bit error probability is minimized. The demodulator uses the maximum-likelihood decision rule to make a *hard decision* on the m -bit code symbol, and the decoder then determines the source symbol from the code symbol. If the code symbols are binary words, Hamming distance is the metric used for decoding, and the code is designed to maximize the minimum Hamming distance.

The concept of *coded modulation* is to accommodate the code redundancy by expanding the signal constellation so that the bandwidth efficiency of the system relative to the uncoded system is unchanged. The channel encoder and signal mapping function of the modulator are designed jointly to maximize the *free distance* (d_{free}) of the code. In the receiver, the decoding and detection operations are also combined so that the decoder operates on *soft decisions*, using the metric for which the code was designed. Section 3.1 illustrates how coded modulation can be used to approach Shannon's capacity bound, and some trade-offs that can be achieved.

Coded modulation was invented in the 1970s by Ungerboeck [1982], who developed schemes, which are known as trellis-coded modulation (TCM), based on convolutional codes. The TCM schemes described by Ungerboeck are optimized for additive white Gaussian noise (AWGN) channels and will be referred to as *Ungerboeck codes*. Ungerboeck encoding involves the convolutional encoding of source symbols followed by the mapping of coded symbols onto an expanded signal set, using a technique called *mapping by set partitioning* to maximize the free Euclidean distance of the code. This process can be represented as a trellis of states with restricted paths, and is described in Section 3.2 with an illustrative example.

The optimum receiver for an Ungerboeck code on an AWGN channel performs soft-decision *maximum-likelihood sequence estimation*, usually using the Viterbi algorithm. Ungerboeck decoding using the Viterbi algorithm is discussed in Section 3.3. It is important to note that the squared Euclidean distance metric, which motivates mapping by set partitioning and is used for decoding, may not be appropriate for channels experiencing

non-AWGN disturbances. The time-dispersive channels investigated in this thesis are of this type. Nevertheless, in the absence of a better metric, the squared Euclidean distance metric is retained.

Coding gain is a measure of the performance of a coded system relative to an uncoded system. TCM can achieve significant *coding gains* on AWGN channels, relative to uncoded systems, without sacrificing data rate or requiring bandwidth expansion. For example, simple 4-state TCM schemes can achieve a 3 dB *asymptotic coding gain*, while more sophisticated TCM schemes have asymptotic coding gains of up to 6 dB. The performance of TCM, decoded using soft-decision Viterbi decoding, is discussed in Section 3.4

We will only study Ungerboeck codes in detail; however, it should be noted that there are other coded modulations that can achieve the same, or higher, coding gains with less complexity. To achieve asymptotic coding gains of 6 dB without using impractical codes, TCM schemes can be designed with multidimensional signal constellations to spread the constellation expansion over more than one signalling interval [Ungerboeck, 1987a; Wei, 1987]. Another variation on Ungerboeck codes is the rotationally invariant codes designed by Wei [1984a; 1984b].

Calderbank and Sloane [1987] have formulated a general class of coded modulations known as *coset codes*, which includes Ungerboeck codes. The signal constellation is regarded as a finite subset of points from an infinite lattice, and set partitioning of the constellation is regarded as a partitioning of the lattice into a sublattice and its cosets. The theory of coset codes has been refined by Forney [1988a; 1988b]. Ungerboeck codes represent some of the best codes in the class of known coset codes. *Multilevel codes* are a generalization of coset codes that can achieve the coding gains of single level coset codes, with reduced complexity decoding [Calderbank, 1989; Pottie and Taylor, 1989a].

3.1 Approaching Shannon's Bound

A natural question to consider is: how closely can coded modulation systems theoretically operate to Shannon's bound? This issue was addressed by Ungerboeck [1982]. Shannon's capacity bound for an AWGN channel is expressed as a bandwidth efficiency bound in (1.5). When a signal constellation is specified and soft-decision decoding is used, Ungerboeck [1982] has shown how to compute the maximum bandwidth efficiency when the signal constellation is used on an AWGN channel. The bandwidth efficiency bound is plotted as a function of signal-to-noise ratio (SNR) in Figure 3.1, together with the maximum bandwidth efficiencies of several signal constellations [Ungerboeck, 1982]. Note that the maximum bandwidth efficiency of the constellations is always less than Shannon's bound.

Shannon predicted that channel coding can be used to achieve an arbitrarily small error probability at any operating point on these curves. For example, a 4-PSK signal constellation can be used in conjunction with channel coding at an SNR of 10 dB to achieve an arbitrarily small error probability at a bandwidth efficiency of close to $\zeta = 2$ bits/s/Hz. Alternatively, an 8-PSK signal constellation used in conjunction with channel coding can achieve the same bandwidth efficiency at a signal-to-noise ratio of 5.9 dB. The result is a 4.1 dB improvement in tolerance to noise, and is within 1.2 dB of the Shannon bound. Note that for $\zeta = 2$ bits/s/Hz, the 16-QAM constellation is only marginally closer to Shannon's bound than the 8-PSK constellation. This observation suggests that there is little to be gained by expanding the constellation by more than a factor of two. The use of signal set expansion and channel coding to improve performance is the essence of coded

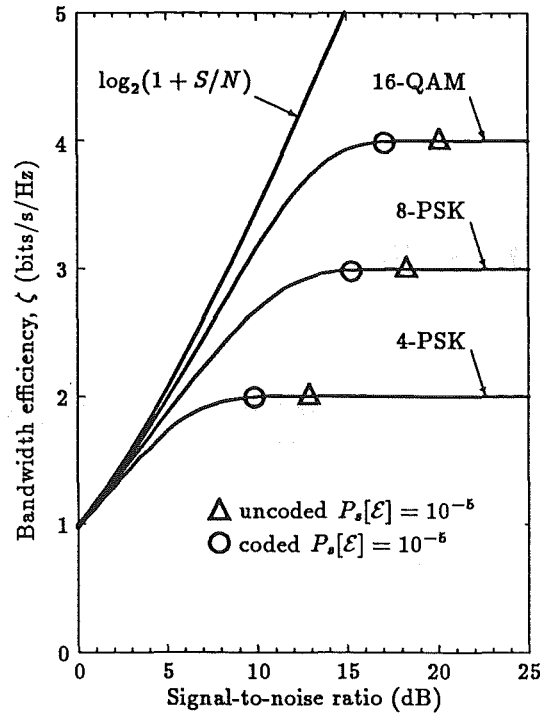


Figure 3.1 Bandwidth efficiency versus signal-to-noise ratio curves.

modulation.

While actual coded modulation schemes have been designed that operate close to the Shannon bound, the *cut-off rate* R_o is regarded by many researchers as a limit to the capacity that can be readily achieved in practice [Forney *et al.*, 1984]. The complexity of sequential decoders is also known to grow rapidly as the code rate is increased above R_o [Wozencraft and Jacobs, 1965]. A detailed discussion of R_o is provided by Wozencraft and Jacobs [1965]. For an AWGN channel

$$R_o = W \log_2(1 + S/2N) \quad (3.1)$$

so R_o is always less than the channel capacity for a given $S/N > 0$. Coded modulations can be designed with asymptotic coding gains of 6 dB; these codes achieve R_o .

If we plot points of equal symbol error probability for the uncoded constellations and typical coded constellations on the bandwidth efficiency curves in Figure 3.1, we can study three *trade-offs* that can be achieved with coded modulation. The Δ markers specify the points where the uncoded constellations have $P_s[\mathcal{E}] = 10^{-5}$. Doubling the number of points in a constellation and using coded modulation allows $P_s[\mathcal{E}]$ to be reduced for the same bandwidth efficiency and SNR. This corresponds to trading complexity (of the code) for improved performance, and can be used in applications that require improved data integrity. Alternatively, coded modulation can be used with an expanded constellation and $P_s[\mathcal{E}] = 10^{-5}$ maintained while the SNR is lowered, as shown by the \circ markers. This corresponds to trading complexity for decreased transmit power, and can be used for power efficient applications; for example, satellite repeaters [Taylor and Chan, 1981]. The third trade-off that can be achieved with coded modulation is to maintain $P_s[\mathcal{E}] = 10^{-5}$ and increase the bandwidth efficiency. This corresponds to trading complexity for bandwidth efficiency, and requires the use of multidimensional constellations to spread the redundancy over more than one signalling interval. Application of this trade-off to voice-band channels has enabled the design of modems for voice-band data transmission

at rates up to 19.2 kbits/s [Pahlavan and Holsinger, 1987].

The trade-offs that can be achieved with coded modulation on time-dispersive channels have not been well studied. This thesis examines the performance improvements that can be achieved by trading the complexity of trellis-coded modulation on channels with time-dispersion and AWGN.

3.2 Ungerboeck Encoding

The structure of an Ungerboeck encoder is shown in Figure 3.2. An m -bit source symbol is transformed into an $m+r$ -bit coded symbol by expanding \tilde{m} ($\leq m$) bits into $\tilde{m}+r$ bits using a binary convolutional encoder of *code rate* $\tilde{m}/(\tilde{m}+r)$, and leaving the remaining $m-\tilde{m}$ bits uncoded. The $2^{\tilde{m}+r}$ coded symbols are then mapped into $2^{\tilde{m}+r}$ signal points on a one-to-one basis. As mentioned in the previous section, there is little to be gained by expanding the signal constellation by more than a factor of two. Therefore, we will henceforth assume that $r=1$, which corresponds to expanding the constellation by a factor of two.

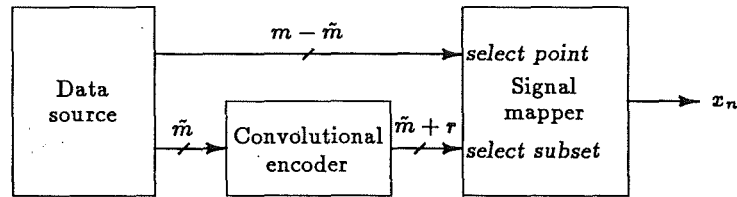


Figure 3.2 Ungerboeck encoder.

The design procedure for a simple Ungerboeck code will be illustrated using a rate 1/2 convolutional encoder and an 8-PSK constellation. Uncoded 4-PSK modulation has the same bandwidth efficiency, and can be regarded as the reference system.

3.2.1 Convolutional Encoding

The rate 1/2 binary convolutional encoder shown in Figure 3.3 consists of $\nu = 2$ memory elements (T-second delays represented by $D = e^{-j\omega T}$) and a modulo two summer. Convolutional encoders are finite-state machines, in this case with four states. The number of output symbols over which a single input has an influence is called the *constraint length* K of the convolutional encoder [Viterbi and Omura, 1979]. The rate 1/2 convolutional encoder has a constraint length of $K = 3$.

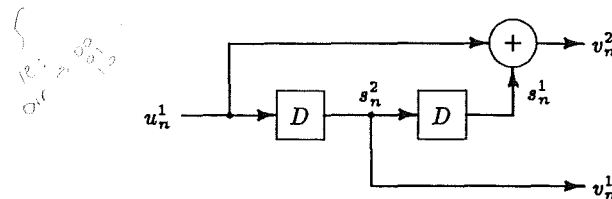


Figure 3.3 Rate 1/2 convolutional encoder.

The operation of a convolutional encoder is completely specified by its trellis diagram, as shown in Figure 3.4 for the encoder in Figure 3.3. The possible states during any symbol period are arranged vertically, and the symbol times are arranged horizontally. For a given present state, the present input symbol determines the next state and the present output symbol. Only a restricted set of paths is possible through the trellis.

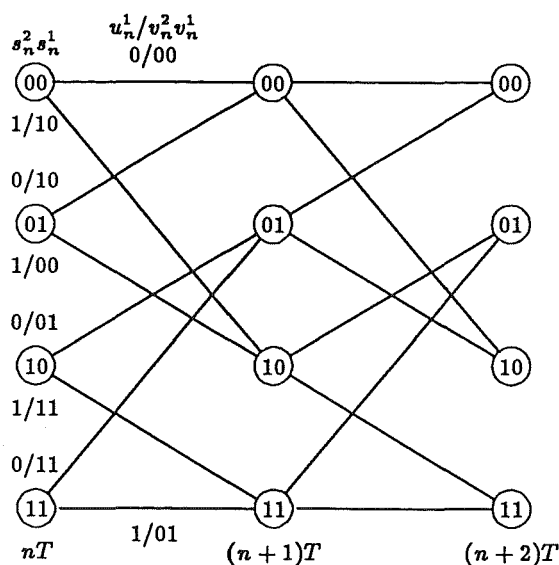


Figure 3.4 Trellis diagram for rate 1/2 convolutional encoder.

An Ungerboeck encoder can be represented by a trellis with state transitions labelled with signal points x_n rather than code symbols v_n . Equivalently, it can be represented by an output function $f(\cdot, \cdot)$ and a next state function $g(\cdot, \cdot)$, such that

$$x_n = f(s_n, u_n) \quad (3.2)$$

and

$$s_{n+1} = g(s_n, u_n) \quad (3.3)$$

where s_n is the present state of the encoder and u_n is the present input to the encoder.

3.2.2 Signal Mapping

The minimum squared distance between two *signal point sequences*, $\mathbf{x} = (x_0, x_1, \dots, x_{N-1})$ and $\hat{\mathbf{x}} = (\hat{x}_0, \hat{x}_1, \dots, \hat{x}_{N-1})$ such that $\hat{\mathbf{x}} \neq \mathbf{x}$, is called the *squared free distance*

$$d_{free}^2 = \min_{\mathbf{x}, \hat{\mathbf{x}} \neq \mathbf{x}, N} \left\{ \sum_{n=0}^{N-1} |x_n - \hat{x}_n|^2 \right\} \quad \text{for } 1 \leq N < \infty \text{ and } \mathbf{x}, \hat{\mathbf{x}} \in \mathcal{X}_N \quad (3.4)$$

Ungerboeck codes are designed to maximize d_{free} using a mapping rule called *mapping by set partitioning*.

Set partitioning involves partitioning a signal set with minimum distance δ_0 between signal points into $2^{\tilde{m}+1}$ subsets with minimum distance $\delta_{\tilde{m}+1}$ between signal points. The procedure is illustrated in Figure 3.5 for an 8-PSK constellation. The signal set is first partitioned into two subsets so that the minimum distance δ_1 between signal points within each subset is maximized. Each of these subsets is then similarly partitioned into two subsets with minimum distance δ_2 , and so on with increasing minimum distances $\delta_0 < \delta_1 < \delta_2 < \dots < \delta_{\tilde{m}+1}$ between the signal points within the subsets.

The $\tilde{m} + 1$ bits from the convolutional encoder determine which one of $2^{\tilde{m}+1}$ subsets is selected, and the $m - \tilde{m}$ uncoded bits determine which of the $2^{m-\tilde{m}}$ signal points within the subset will be selected. The effect of the $m - \tilde{m}$ uncoded bits is to introduce *parallel transitions* into the trellis.

The trellis diagram for our example of coded 8-PSK is shown in Figure 3.6. Three rules were formulated by Ungerboeck [1982] for assigning signal points to trellis state transitions:

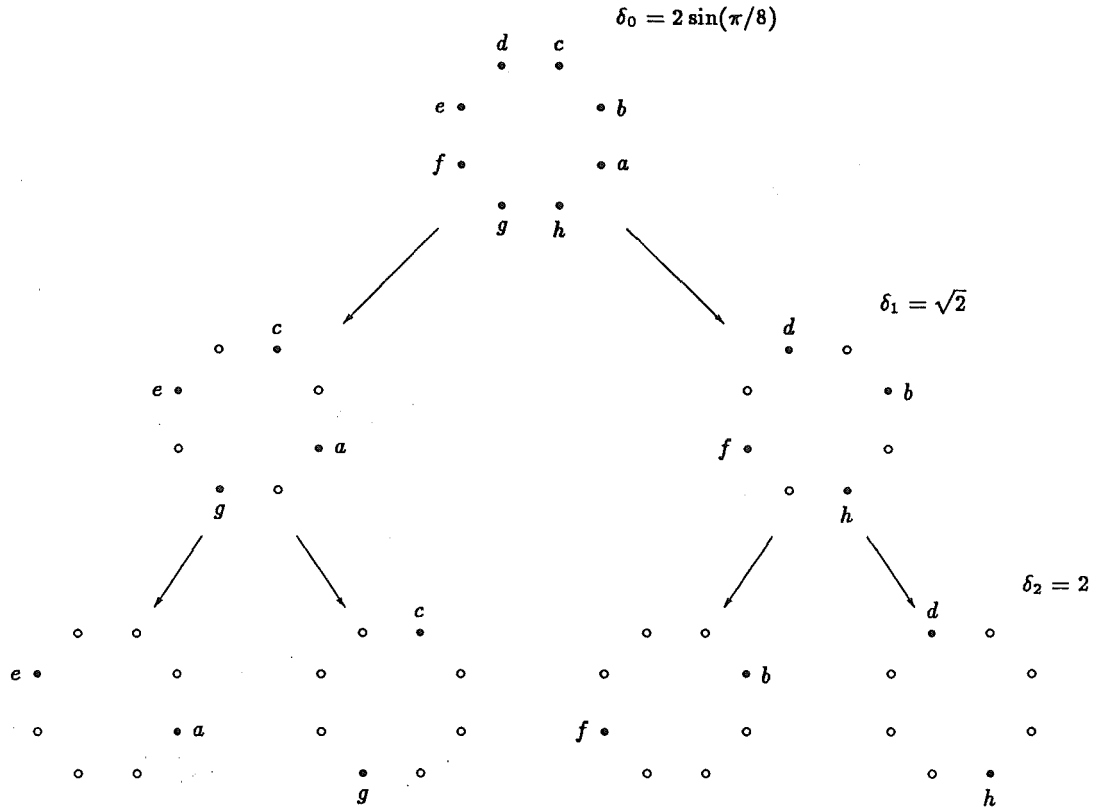


Figure 3.5 Set partitioning of 8-PSK.

1. Parallel transitions are associated with signal points with maximum distance $\delta_{\tilde{m}+1}$ between them. In our example, the parallel transitions are associated with signal points from the subsets (a, e) , (c, g) , (b, f) , or (d, h) .
2. Transitions originating from or merging into one state are associated with signal points with the next maximum possible distance between them. In our example, transitions originating from or merging into the same state are associated with signal points from the subsets (a, c, e, g) or (b, d, f, h) .
3. All signal points are used in the trellis diagram with equal frequency.

The free distance of the coded 8-PSK can be determined from the trellis diagram in Figure 3.6. The minimum squared Euclidean distance between two signal point sequences that diverge at one state and merge at another state, after more than one transition, is $\delta_1^2 + \delta_0^2 + \delta_1^2 (= \delta_2^2 + \delta_0^2)$; the two paths with signal point sequences aaa and cbc (shown in bold) have this distance between them. The squared distance between these paths is greater than the distance of $\delta_2 = 2$ between signal points assigned to parallel transitions. As defined by (3.4), the minimum distance between two different paths is the free distance d_{free} ; in this case $d_{free} = 2$. This free distance is an improvement of 3 dB over the minimum distance of $\sqrt{2}$ between the signal points of uncoded 4-PSK. To increase d_{free} when it is associated with parallel transitions, without introducing more redundancy or increasing the average energy in the signal constellation, more source bits must be encoded because increasing the number of code states does not increase the distance between parallel transitions.

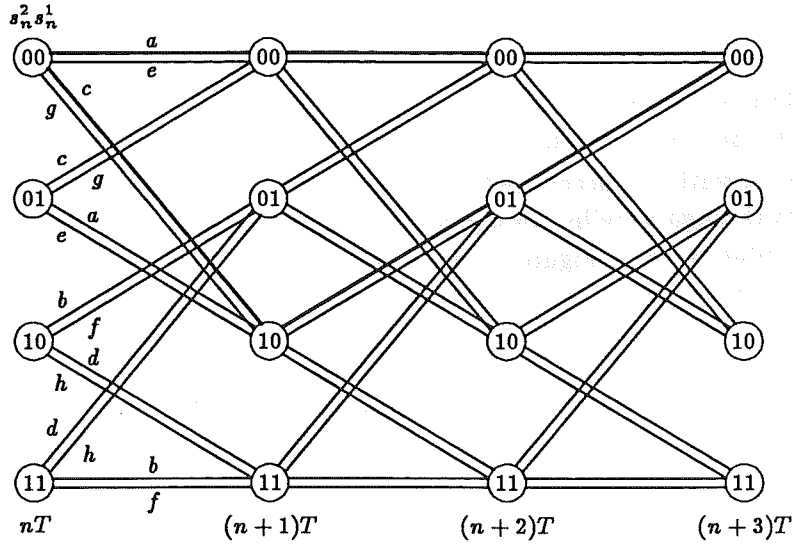


Figure 3.6 Trellis diagram for coded 8-PSK.

3.3 Ungerboeck Decoding

The *hard* receiver decisions, made prior to decoding in a conventional receiver, result in an irreversible loss of information because information on *how close* the received signal was to the selected signal point is discarded. To improve the performance of a receiver when the transmitted signal points are mutually dependent, the decoder can be designed to operate on unquantized received signal samples. Such decoding is called *soft-decision decoding*, and a squared Euclidean distance metric is appropriate for decoding and code design on an AWGN channel. Soft decisions provide about 2 dB improvement over the performance possible with hard decisions [Heller and Jacobs, 1971].

Consider the received signal samples $r_n = x_n + \eta_n$, where at time n , x_n is the signal point generated by the modulator, and η_n is a sample of an AWGN process. The decision rule of the optimum receiver for an Ungerboeck code determines, from the set \mathcal{X}_N of all signal point sequences, the sequence $\hat{\mathbf{x}} = (\hat{x}_0, \hat{x}_1, \dots, \hat{x}_{N-1})$ with the minimum squared Euclidean distance d_{min}^2 from $\mathbf{r} = (r_0, r_1, \dots, r_{N-1})$. This is the sequence $\hat{\mathbf{x}}$ that satisfies

$$d_{min}^2 = \min_{\hat{\mathbf{x}}} \left\{ \sum_{n=0}^{N-1} |r_n - \hat{x}_n|^2 \right\} \quad \text{for } \hat{\mathbf{x}} \in \mathcal{X}_N \quad (3.5)$$

Provided the transmitted signal point sequence $\mathbf{x} \in \mathcal{X}_N$ has been generated according to the rules of a finite-state machine, the *Viterbi algorithm* can be used to determine the signal point sequence $\hat{\mathbf{x}}$ closest to the received signal sequence \mathbf{r} .

The complexity of the Viterbi algorithm grows exponentially with the number of bits encoded (\tilde{m}) and the number of memory elements in the encoder (ν). In practice, the code rate $\tilde{m}/(\tilde{m} + 1)$ and the number of code states are limited by the data rate at which the Viterbi algorithm must operate. Although we will use the Viterbi algorithm in our work, it should be noted that there are a variety of reduced complexity, but suboptimum, techniques for decoding convolutional codes. These techniques can also be applied to decode TCM, as described by Pottie and Taylor [1989b].

3.3.1 Viterbi Decoding

The Viterbi algorithm was developed by Viterbi [1967] as an 'asymptotically optimum' technique for decoding convolutional codes. With soft-decision Viterbi decoding, the most likely signal point sequence is determined directly from the unquantized received signal.

The Viterbi algorithm [Forney, 1973; Hayes, 1975] is an efficient technique to find the most likely path through a trellis, for a given received signal sequence. Viterbi decoding of coded 8-PSK is illustrated in Figure 3.7 for a received signal sequence $aabaaaa$, which was transmitted as $aaaaaaa$. Note that the received signal sequence consists of signal points because hard-decision Viterbi decoding is used to simplify the example. The algorithm requires a distance metric and path history to be stored for each state. Assume that the decoder knows the correct trellis state at time zero (this assumption is not necessary in practice, but simplifies the explanation). At time zero, all states are assigned a metric of zero. For each state at time T , the most likely transition to the state is determined as the one associated with the signal point that is closest, in terms of Euclidean distance, to the received signal, and its squared distance is added to the state metric. The signal point associated with the most likely transition is stored in the path history. At time $(n+1)T$, the algorithm extends the paths from time nT by choosing the best path to each new state as a *survivor* and forgetting all other paths that cannot be extended as best paths to the new states.

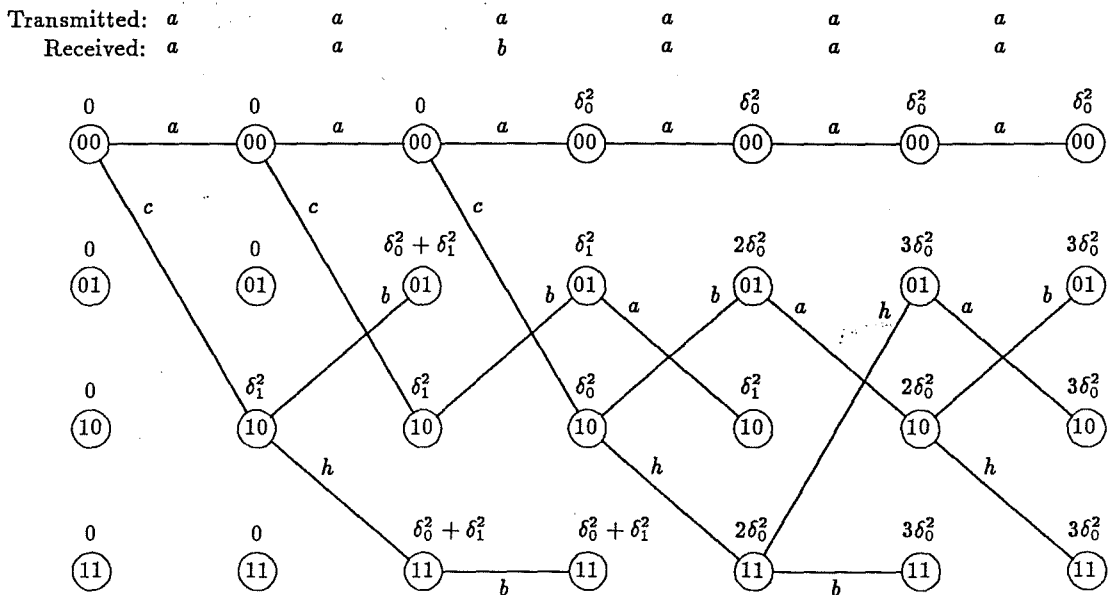


Figure 3.7 Trellis diagram for Viterbi decoding of the hard-decision received sequence $aabaaaa$.

The surviving paths at time nT tend to merge into the same single *history path* at time $(n - \Delta)T$. With a sufficiently large decoding delay Δ (usually four or five code constraint lengths [Heller and Jacobs, 1971]), there is a high probability that all paths will have merged, and the signal point associated with the merged transitions can be taken as the most likely signal point. Thus, in practice, the decoder can make decisions after a finite decoding delay, even though an infinite delay is theoretically required. The maximum-likelihood signal point sequence is then mapped back to an uncoded symbol sequence. This mapping is normally performed within the Viterbi algorithm itself.

If the surviving paths do not merge into the same state at the decoding depth, the decoder must make a *forced decision*, using an arbitrary rule such as choosing the path

with the lowest metric, choosing the path through the most common state at the decoding depth, or choosing any candidate path. Alternatively, the decoder can retain a record of the tie and use it to flag an uncorrectable codeword segment.

When implementing the Viterbi algorithm with soft decisions, true soft decisions cannot be used because they require infinite precision. A resolution of eight quantization levels between signal points incurs a loss of less than 0.25 dB relative to absolute precision soft-decision decoding [Heller and Jacobs, 1971], and only requires three bits per dimension in addition to the bits required for hard decisions.

In practice, soft-decision Viterbi decoding for TCM involves two steps. The decoder first determines the closest signal point to the received signal within each subset. This is known as *subset decoding*, and can be performed with look-up tables. The signal point selected from each subset and its squared distance metric are then used in the Viterbi algorithm to determine the most likely sequence of subset decoded signal points. This procedure is equivalent to, but simpler than, implementing the Viterbi algorithm with the full signal set.

3.4 Performance of Trellis-Coded Modulation

The most probable errors made by a soft-decision Viterbi decoder occur between transmitted and decoded signal point sequences \mathbf{x} and $\hat{\mathbf{x}}$ that are closest in terms of Euclidean distance. An *error event* occurs in the receiver when the Viterbi algorithm makes a wrong decision, causing the decoded symbol sequence to diverge and later merge with the transmitted symbol sequence. An error event is illustrated in Figure 3.8 for a transmitted signal point sequence $aaaaaa$ and decoded signal point sequence $aacbca$.

The error-event probability on an AWGN channel is lower bounded by [Forney, 1973]

$$P_e[\mathcal{E}] \geq N(d_{free}) Q\left(\frac{d_{free}}{2\sigma_\eta}\right) \quad (3.6)$$

where σ_η^2 is the variance of the white Gaussian noise samples, $N(d_{free})$ is the average number of nearest-neighbour signal point sequences at the free distance d_{free} , and $Q(\cdot)$ is the Gaussian integral function defined in (2.29). At high signal-to-noise ratios, the dominant error event is the one at the free distance, and the lower bound asymptotically approaches

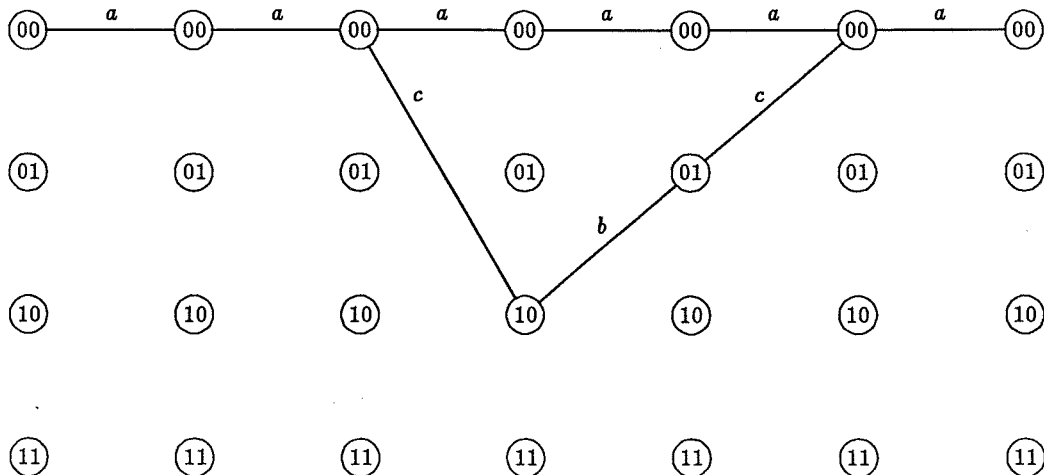


Figure 3.8 An error event.

the actual error-event probability. A lower bound on the symbol error probability $P_s[\mathcal{E}]$ can be obtained by multiplying the lower bound on the error-event probability by the average number of symbols per d_{free} error event.

Viterbi [1971] developed a union bound on the error-event probability of the Viterbi algorithm used to decode a convolutional code. This bound takes advantage of the algebraic structure of convolutional codes, which allows the union bound to be formed by only considering a single transmitted symbol sequence (reference path). It can be evaluated using the transfer function of the code, with a cost proportional to the number of code states. Biglieri [1984] extended this technique to trellis-coded modulation by considering all possible transmitted symbol sequences. However, the cost is proportional to the square of the number of code states. Zehavi and Wolf [1987] found that, for certain classes of TCM, the cost of computing the union bound could be reduced to being proportional to the number of code states. They used a modified code transfer function to evaluate the union bound.

The performance of coded systems relative to equivalent uncoded systems is often measured by the *coding gain* in decibels on an AWGN channel. Coding gain is defined at a specific bit, symbol, or event error probability, and is the difference between the SNR of the uncoded system and the SNR of the coded system when they operate at the same error probability. The *asymptotic coding gain* approached at high signal-to-noise ratios by a TCM system relative to an equivalent uncoded system is given in decibels by [Ungerboeck, 1987b]

$$G_{c,u} = 10 \log_{10}(d_{free,c}^2/d_{free,u}^2) - 10 \log_{10}(E_{s,c}/E_{s,u}) \quad (3.7)$$

where $d_{free,c}^2$ and $d_{free,u}^2$ are the squared free distances, and $E_{s,c}$ and $E_{s,u}$ denote the average energy in the signal constellations for the coded and uncoded systems respectively. The first term is the gain due to increased free distance, and the second term is the penalty due to increased average energy in the signal constellation. Note that asymptotic coding gain is defined in terms of error-event probabilities. The asymptotic coding gain of coded 8-PSK relative to uncoded 4-PSK is 3 dB.

3.5 Conclusion

In this chapter an introduction to TCM has been provided and an attempt has been made to show why TCM is a powerful coding technique. The codes we study in this thesis will be Ungerboeck codes, although one of the codes is nonlinear and uses Wei's rules for rotational invariance. We study a range of signal constellations from 4-PSK to 1024-QAM. Viterbi decoding is used for all the codes, and reduced complexity decoding techniques are not studied. We examine the performance improvements due to TCM on time-dispersive channels. In particular, we are interested in the performance improvements that TCM can offer with the residual ISI that exists after non-ideal equalization.

Chapter 4

Digital Microwave Radio Systems

Terrestrial microwave radio systems were first developed to carry long distance voice and video signals. These systems used *analog* frequency modulation and, therefore, required additional terminal equipment to carry data. In response to the increasing use of the public switched network for data transmission, *digital* microwave radio (DMR) systems were deployed in the early 1970s. Digital systems achieve higher quality transmission than their analog counterparts and, because they use repeaters for signal regeneration, their performance is essentially independent of distance.

The first DMR systems used 4-PSK modulation and achieved bandwidth efficiencies of less than 2 bits/s/Hz. The basic components of such systems will be discussed in Section 4.1. The 1980s saw the deployment of 16-QAM and then 64-QAM systems. During the 1990s we will see the introduction of commercial 256-QAM systems, with bandwidth efficiencies of about 7 bits/s/Hz, and even higher level systems that are currently in the planning stages [Komaki *et al.*, 1990].

Optical fibre systems present a major technological challenge to the future of DMR systems because they already offer greater capacity and promise almost unlimited capacity in the future. In spite of this, the different characteristics of optical fibre systems and DMR systems are complementary. Optical fibre systems can provide very high capacity (in excess of 1 Gbits/s) at a lower cost per unit bandwidth than DMR systems. In contrast, low to high capacity (2 to 400 Mbits/s) systems can be provided more cheaply with DMR technology. DMR systems are still the backbone of many telecommunication networks, and the investment in antenna towers and buildings often allows new systems to be installed even more cost effectively. A further advantage of DMR systems is the network security they provide; except for the cables to the nearest switching offices, there are no cables in the ground to be damaged. To ensure the future of DMR, it must be integrated with the rest of the digital network.

The high data rates carried by many DMR systems challenges the realization of hardware—particularly equalizers and decoders. Signal processing technology that is easily implemented for voiceband modems ($R < 20$ kbits/s) is often unable to be implemented at these data rates.

The most troublesome disturbance on the DMR channel is not noise (signal-to-noise ratios of 60 dB are readily achieved in the absence of flat fading), but rather *multipath fading*, which can introduce *frequency selective* notches or slopes that cause severe intersymbol interference (ISI) in the receiver. Although frequency selective fading is not a problem with 4-PSK systems, it is a source of serious performance degradation for systems with higher bandwidth efficiency. Section 4.2 discusses multipath fading and channel

modelling.

DMR systems are also subject to adjacent channel interference, co-channel interference, rain attenuation, and ducting [Townsend, 1988]. These tend to be less problematic than the intersymbol interference induced by frequency selective fading and will not be further considered.

The radio spectrum is a costly resource because it is in high demand; therefore, the bandwidth efficiency of DMR systems must continue to increase in the interests of economy. The main technique for improving bandwidth efficiency is to increase the number of points in the signal constellation. If the Euclidean distance between signal points is to be kept constant so the immunity to Gaussian noise is unchanged, the cost of increasing bandwidth efficiency is increased transmit power (power efficiency is less important than bandwidth efficiency for DMR systems). There is an additional cost on frequency selective channels, due to the larger number of signal points in the constellation, which results in more severe ISI and higher error rates. To maintain a satisfactory error rate, countermeasures must be introduced to combat the ISI induced by the fading. Countermeasures for multipath fading are discussed in Section 4.3.

A technique that effectively doubles bandwidth efficiency is to transmit independent information on vertically and horizontally polarized waves. Although, in terms of information theory, the vertical and horizontal polarizations should be considered as two channels of equal bandwidth, so the bandwidth efficiency is not changed. For constellations with more than sixteen signal points, the isolation provided by typical antenna cross polarization discriminations is not sufficient, and adaptive cross polarization interference cancellers are required. These cancellers must be able to cope with frequency selective fading because cross polarization interference tends to be most severe during periods of multipath fading. In this thesis, cross polarization and the effects of cancellers are not studied.

It is important to be able to measure the performance of operating digital systems and to estimate the performance of proposed systems. The dynamic fading nature of the DMR channel makes the definition of a measure of performance difficult. The most popular measure is link outage, which is discussed in Section 4.4.

This chapter reviews the aspects of DMR that are significant to the main thrust of this thesis. A more detailed review is provided by the tutorial series on DMR that appeared in IEEE Communications Magazine from August 1986 to February 1987 [Taylor and Hartmann, 1986; Noguchi *et al.*, 1986; Rummler *et al.*, 1986; Chamberlain *et al.*, 1986; Greenstein and Shafi, 1987; Yamamoto, 1987]. Townsend [1988] and Ivanek [1989] provide comprehensive references on DMR systems.

4.1 Modulation and Demodulation

The term modulation is often used broadly to refer to the signal mapping, pulse shaping, and carrier modulation functions of a transmitter. Similarly, the term demodulation often refers to the detection, matched filtering, and baseband recovery functions of a receiver. The functions of a basic *modem* (modulator/demodulator), with no countermeasures for multipath fading, are discussed in this section.

4.1.1 Modulation

Figure 4.1 shows the schematic of a basic modulator. The serial bit sequence, input to the modulator, is converted into the *in-phase* (I) and *quadrature* (Q) binary symbol sequences.

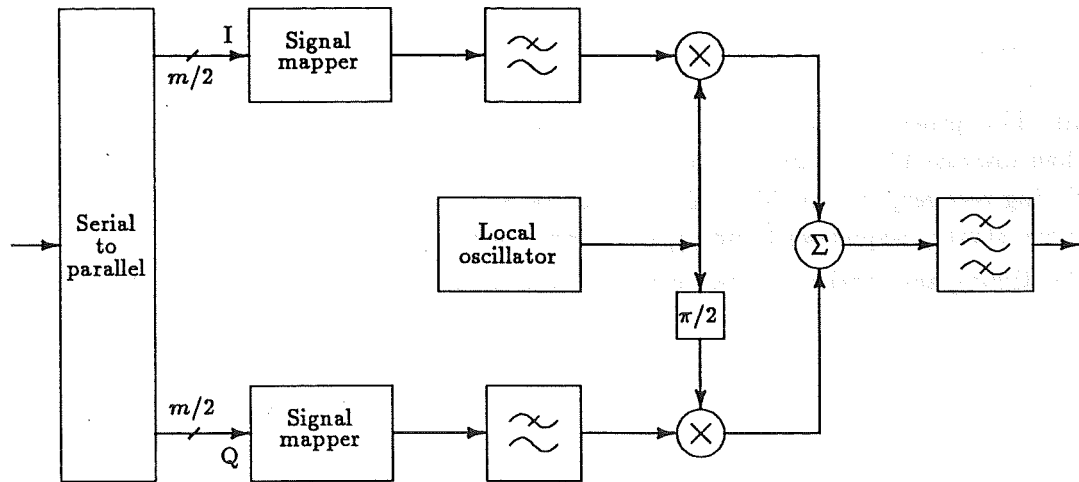


Figure 4.1 Schematic of a basic modulator.

These I and Q symbol sequences are then mapped into multilevel analog pulse sequences, which are lowpass filtered. Two sinusoidal carriers—derived from the *local oscillator*—at the same frequency but in phase quadrature are modulated by these pulse sequences. The modulated signals are then added and bandpass filtered. In practice, signals are usually modulated onto an *intermediate frequency* (IF) carrier to simplify hardware realization, and the IF signal is translated to a *radio frequency* (RF).

The I and Q carriers are orthogonal, and can be used to transmit any two-dimensional signal constellation (or higher dimensionality if time orthogonality is also used). However, PSK and QAM are the most widely used constellations for DMR systems.

4.1.2 Demodulation

Figure 4.2 shows the schematic of a basic demodulator. Initially the received signal is bandpass filtered to select the desired signal and provide noise rejection. The sinusoidal

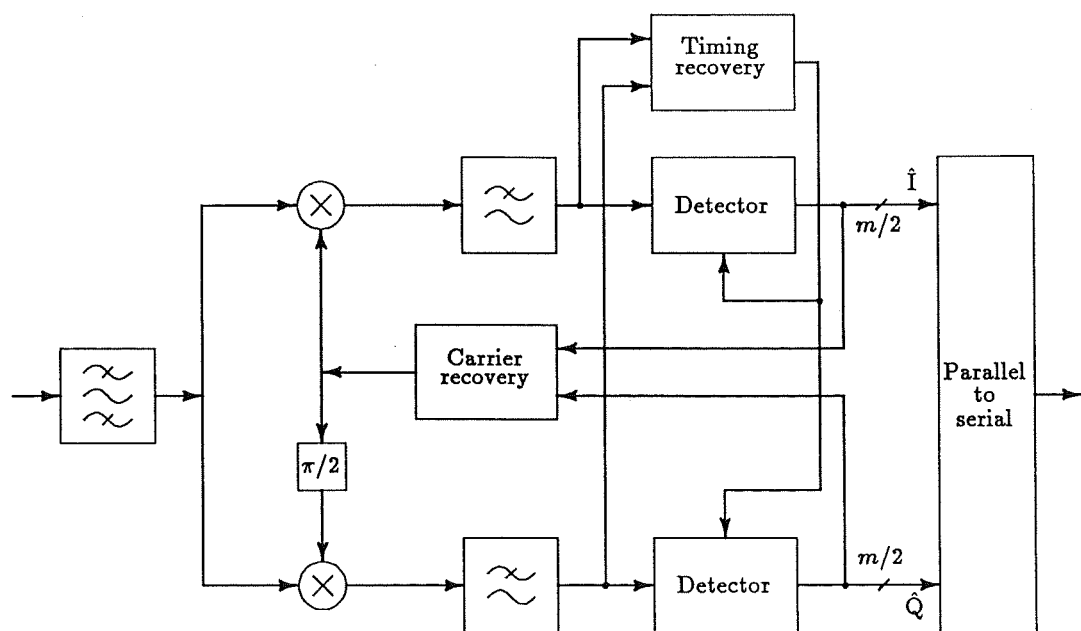
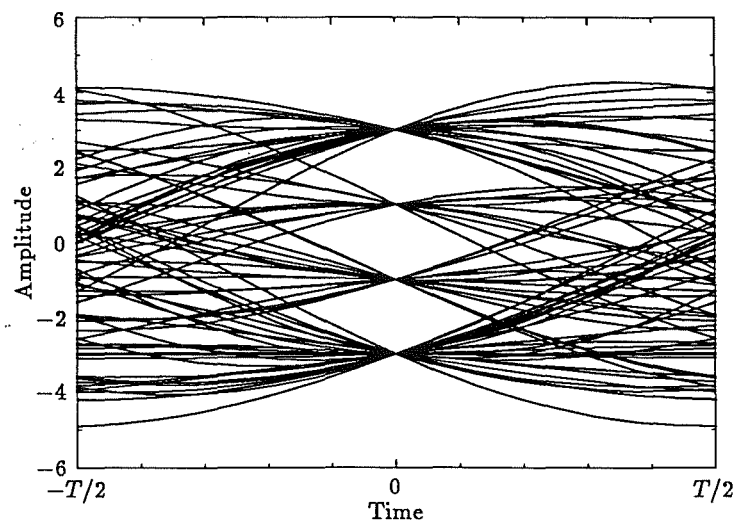


Figure 4.2 Schematic of a basic demodulator.

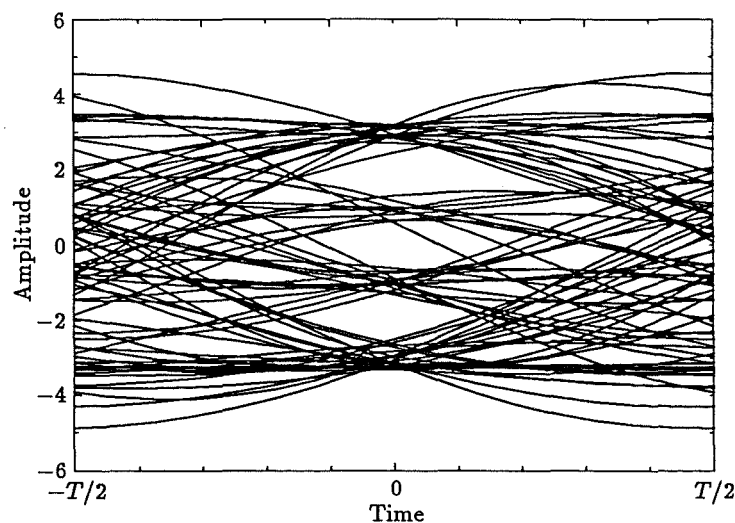
carrier is removed from the received signal by multiplying it by quadrature sinusoids, which are phase and frequency locked to the transmitter local oscillator by the carrier recovery circuit. This process is known as *coherent demodulation*. The resultant baseband signals are then lowpass filtered and sampled once per symbol, at an optimum time specified by the timing recovery circuit. The sampled analog signal is converted by a detection circuit to binary symbol sequences \hat{I} and \hat{Q} , which are then converted to a serial bit sequence.

The timing and carrier recovery circuits must extract their references from the received signal, so should be robust to the distortions introduced by multipath fading. A number of techniques can be used for carrier and timing recovery; these are reviewed by Franks [1980].

The effects of timing errors and frequency selective fading can be seen by examining *eye diagrams* [Proakis, 1989] like those in Figure 4.3 for 16-QAM. Optimum decision thresholds are located at the vertical midpoint of the eyes, and the optimum sampling time is located at the horizontal midpoint of the eyes. Figure 4.3a is the eye diagram without ISI at the optimum sampling time. In this situation, the receiver has maximum



(a) Without ISI



(b) With ISI

Figure 4.3 Eye diagrams of the I channel for 16-QAM.

tolerance to noise. Any timing error causes the receiver to sample at a time when the eye is not fully open, resulting in ISI. The effect of ISI caused by frequency selective fading is shown in Figure 4.3b. The ISI reduces the eye opening height even at the optimum sampling time, so the receiver becomes more susceptible to errors induced by noise. If the ISI is sufficiently severe, the eyes will close, and the ISI alone will cause decision errors.

Carrier phase errors have the effect of rotating the signal constellation and introducing interference between the I and Q signals.

4.1.3 Pulse Shaping

To avoid ISI, it is necessary that the combined response of the square analog pulses, transmitter filters, receiver filters, and channel satisfy Nyquist's criterion for zero ISI (see Section 1.1.1). In practice, the channel impulse response may change with time and is often unknown. There are two solutions to this problem; either the receiver can be designed to adaptively cancel the effect of the channel or the channel can be assumed to be perfect (the DMR channel is nearly perfect most of the time). Adaptive receivers can achieve better performance than fixed receivers, but at the cost of increased complexity.

The most popular Nyquist pulse shapes are from the family of raised-cosine pulses with impulse response

$$h(t) = \frac{\sin(\pi t/T)}{\pi t/T} \frac{\cos(\pi \alpha t/T)}{1 - (2\alpha t/T)^2} \quad (4.1)$$

where $0 \leq \alpha \leq 1$ is called the *roll-off factor*. The frequency response of these pulses is

$$H(f) = \begin{cases} T & \text{if } |f| < (1 - \alpha)/2T \\ T \cos^2 \left(\frac{2\pi T |f| - \pi(1 - \alpha)}{4\alpha} \right) & \text{if } (1 - \alpha)/2T \leq |f| \leq (1 + \alpha)/2T \\ 0 & \text{if } |f| > (1 + \alpha)/2T \end{cases} \quad (4.2)$$

Compromises must be made when choosing α . If α is near one, the bandwidth efficiency will be half that theoretically attainable. If α is near zero, manufacturing the filters will be more difficult and costly, and the pulse sidelobes will be higher, making the received signal more susceptible to time-dispersion. The most common values for α are $0.3 \leq \alpha \leq 0.5$.

Another factor to be considered is the partitioning of the filter response between the transmitter and the receiver. The usual allocation of filtering is as follows. The transmit lowpass filter is designed to implement the combination of a root Nyquist filter $\sqrt{H(f)}$ and an inverse sinc filter $2\pi fT/(\sin 2\pi fT)$, to convert the square pulses, of duration T seconds, to impulses. The transmit bandpass filter is designed for sideband rejection. The receive bandpass filter is designed for noise rejection and signal selection, while the receive LPF is designed as a root Nyquist filter $\sqrt{H(f)^*}$ (* denotes complex conjugation) and performs sideband rejection. The receive filter frequency response is the complex conjugate of the transmit filter frequency response, so that the impulse response of the receive filter is a time reversal of the transmit filter. This produces a matched filter at the receiver in the absence of fading.

Nyquist filters can be implemented at a translated frequency, but their realization is more difficult than at baseband because of the additional requirement of conjugate symmetry. Analog or digital implementations can be used at baseband frequencies, but analog implementations must be used at IF because a digital implementation would require sampling at twice the highest frequency in the IF signal. Thus a sampling rate in excess of twice the IF would be necessary.

4.2 Multipath Fading

DMR systems usually experience normal—unfaded—propagation conditions for a large proportion of time. An AWGN channel essentially exists under these conditions, and the design margins to guard against outage during periods of fading allow bit error rates of less than 10^{-10} to be achieved.

Under normal propagation conditions, the refractive index of the atmosphere decreases with height, and the propagation path from transmitter to receiver curves approximately with the surface of the earth as illustrated in Figure 4.4a. Anomalous atmospheric conditions can lead to the formation of inhomogeneous temperature and humidity profiles. The resulting layered atmosphere causes sharp changes in the vertical refractive index gradients, giving rise to multiple propagation paths with different attenuations and delays. Thus the transmitted signal arrives at the receiver via a number of paths as illustrated in Figure 4.4b. These multiple paths cause the transmitted signal to be distorted at the receiver. Multipath fading occurs mainly in temperate climates on hot humid summer evenings with no clouds and little wind.

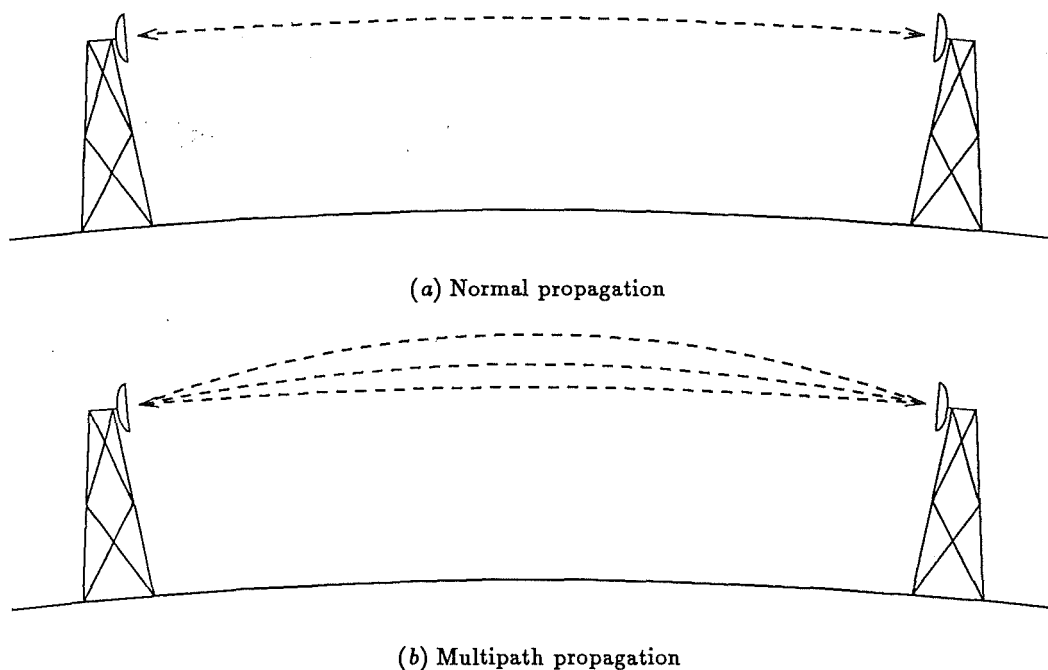


Figure 4.4 Propagation conditions.

To analyze the effect of multipath fading on system performance, it is desirable to develop a model of the channel under this condition. We will consider *channel models*, which describe the frequency response of the channel, rather than *atmospheric models*, which describe the physical propagation. Siller [1984] and Rummler *et al.* [1986] provide excellent reviews of multipath fading and models for the DMR channel; they also give comprehensive lists of references.

4.2.1 Channel Models

The most popular channel models derive from the *simplified three-path model*, which has a transfer function

$$H(f) = a \left[1 - be^{-j2\pi(f-f_0)\tau} \right] \quad (4.3)$$

where a is the flat fade term, b is the depth of the notch that occurs at frequency f_o (notch frequency), and τ is the relative delay between paths. In the time domain, the simplified three-path model becomes

$$h(t) = a [\delta(t) - be^{j2\pi f_o \tau} \delta(t - \tau)] \quad (4.4)$$

This model is referred to as a three-path model because it considers a direct unfaded path, a second path close enough in delay so the combined response of these two paths has a constant magnitude a over the channel bandwidth, and a third path at relative delay τ , which provides the frequency shaping of $H(f)$. Whenever $b \neq 0$, the fading introduces notches in frequency, and the fading is frequency selective if $H(f)$ is distorted over the channel bandwidth. Such fades are time-dispersive and cause ISI.

The zeros in $H(f)$ can be examined by expressing (4.3) in the Laplace domain to get

$$H(s) = a [1 - be^{j2\pi f_o \tau} e^{-s\tau}] \quad (4.5)$$

where $s = \sigma + j2\pi f$. The zeros in this response occur at

$$s = \ln(b)/\tau + j2\pi (f_o \pm k/\tau) \quad k = 0, 1, \dots \quad (4.6)$$

as illustrated by the zero-maps in Figure 4.5. If $\tau > 0$ and $b < 1$, or $\tau < 0$ and $b > 1$, the zeros lie in the left-half plane as shown in Figure 4.5a, and the response is minimum phase. If $\tau < 0$ and $b < 1$, or $\tau > 0$ and $b > 1$, the zeros lie in the right-half plane as shown in Figure 4.5b, and the response is non-minimum phase. The closer b is to unity,

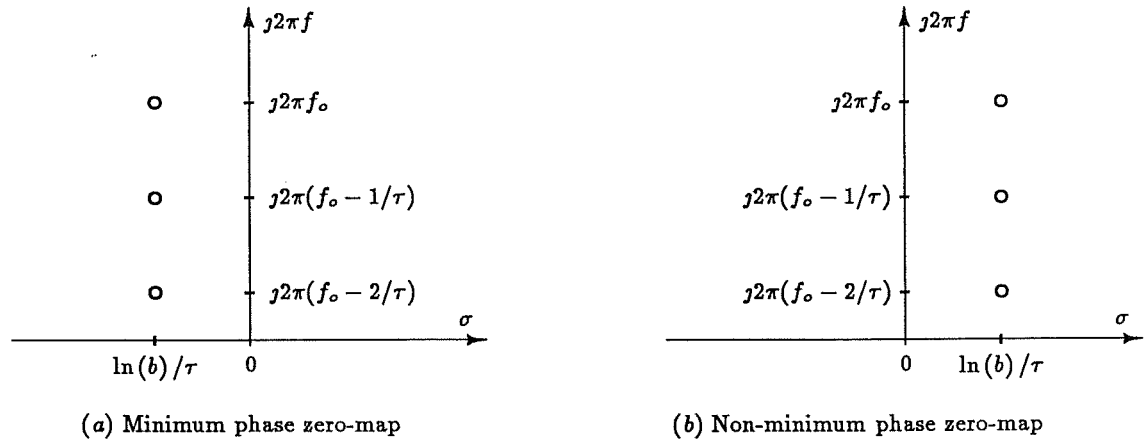


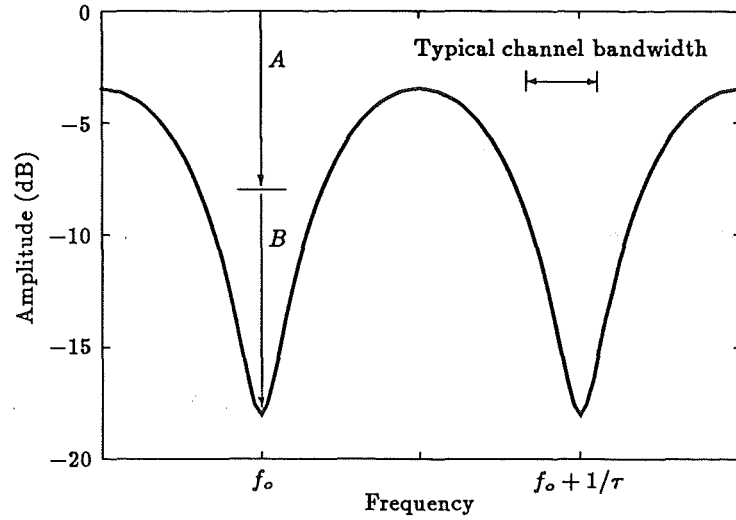
Figure 4.5 Zero-maps for the simplified three-path model.

the closer the zeros move to the $j2\pi f$ axis, and the deeper are the notches in $H(f)$.

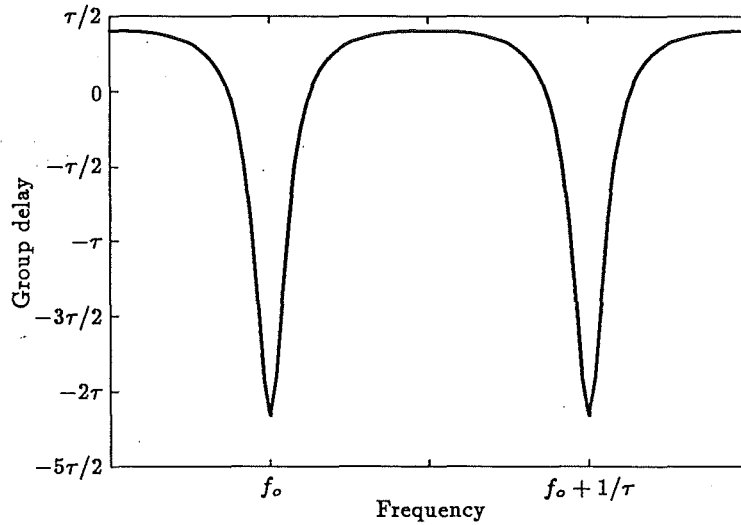
In general, minimum phase fading occurs more frequently than non-minimum phase fading for shallow fading events, and minimum and non-minimum phase fading occur with equal frequency for severe fading events [Rummler *et al.*, 1986]. Some countermeasures for multipath fading are not equally effective against minimum and non-minimum phase fading. For example, a *decision-feedback equalizer* is more effective against minimum phase fading than non-minimum phase fading (see Section 4.3.1).

The distortion caused by fading can be specified in terms of amplitude distortion and group delay distortion. Group delay is defined in terms of the phase response $\angle H(f)$ of $H(f)$ as

$$D(f) \triangleq -\frac{\partial \angle H(f)}{2\pi \partial f} \quad (4.7)$$



(a) Amplitude response



(b) Group delay response

Figure 4.6 Typical amplitude and group delay responses for $H(f)$.

A typical amplitude response of $H(f)$ is shown in Figure 4.6a, and a typical group delay response of $H(f)$ is shown in Figure 4.6b. Note that only a single notch can occur within a typical channel bandwidth. Parameters a and b are expressed in decibels by

$$\begin{aligned}
 A &\triangleq -20 \log_{10}(a) \\
 B &\triangleq -20 \log_{10}(1 - b) \quad \text{if } 0 \leq b \leq 1 \\
 &\triangleq -20 \log_{10}(1 - 1/b) \quad \text{if } b > 1
 \end{aligned} \tag{4.8}$$

To avoid amplitude distortion, the amplitude response should be constant across the channel bandwidth. Group delay distortion can be avoided if the group delay is constant across the channel bandwidth.

The three-path channel modelling function has been shown to provide a good fit to measured responses over a channel of bandwidth W when $\tau < 1/6W$ [Rummler, 1979], which is the case for the DMR channel. Most channel defects can be described as either

attenuation slopes (caused by out-of-band notches) or single notches. In practice, the three-path model actually has too many degrees of freedom for these channel conditions. Within measurement errors, a unique set of model parameters cannot be determined for a given measured channel response. To overcome this problem, Rummler *et al.* [1986] used a fixed value of $\tau = 6.3$ ns. Thus we have the *Rummler model* for the DMR channel.

When the flat fading component is not of interest, the simplified three-path model can be reduced to a *two-path model* with frequency response

$$H(f) = 1 - be^{-j2\pi(f-f_o)\tau} \quad (4.9)$$

This is a purely time-dispersive fading model because it only considers a direct ray and a dominant interfering ray. It is useful when dispersive fading dominates over flat fading.

4.3 Countermeasures for Multipath Fading

The use of constellations with more than sixteen signal points has led to the development of countermeasures to combat the undesirable effects of fading.

Flat fading causes the received power level to vary; therefore, an automatic gain control (AGC) is normally used to maintain constant power in the receiver. The AGC, however, does not improve the degraded signal-to-noise ratio (SNR) that results from flat fading. This degraded SNR can be improved with diversity or error-control coding.

To combat dispersive fading, we can use equalization, diversity, and error-control coding. Adaptive equalization is now included, to combat dispersive fading, in all DMR systems with constellations of sixteen or more signal points. Space and frequency diversity techniques are also popular as countermeasures to dispersive fading. Error-control coding techniques offer additional performance improvements by cleaning up the errors due to residual ISI after other countermeasures have been applied. This section discusses adaptive equalization, diversity, and error-control coding.

4.3.1 Adaptive Equalization

Adaptive equalization can be performed in either the time or frequency domains. Time-domain equalization is the more natural form of equalization because the equalizer acts directly to reduce ISI. It is also the most effective equalization technique.

Frequency-domain equalization is implemented at IF and adapted by monitoring the power spectrum at two or three frequencies, using a set of narrowband filters and comparing the measured powers with each other or with pre-determined undistorted levels maintained by the AGC. The first frequency-domain equalizers to be developed were linear amplitude-slope equalizers, which attempt to remove the spectral slope across the system bandwidth that can result from out-of-band notches. These equalizers provide good performance improvement for out-of-band notches, but provide little compensation for in-band notches.

To compensate for in-band notches, notch equalizers were developed. Notch equalizers are resonant filters that track and attempt to remove a spectral notch within the system bandwidth. These equalizers have concave (\cap) group delay characteristics and can achieve significant reductions in distortion for minimum phase fading, which has convex (\cup) group delay. However, the group delay distortion is doubled for non-minimum phase fading, and performance is similar to that of a slope equalizer [Chamberlain *et al.*, 1986].

The simplest time-domain equalizer is an adaptive transversal filter (tapped delay-line) as shown in Figure 4.7. The output of the equalizer $y(nT)$ is produced by summing

linearly weighted samples of the received signal $r(nT)$. The equalizer tap weights are w_i , $-k_1 \leq i \leq k_2$, where k_1 is the number of weights to combat precursor distortion and k_2 is the number of weights to combat postcursor distortion. This particular equalizer is known as a *synchronously-spaced* equalizer (SSE) because there is one tap per symbol (T -spacing). DMR systems typically use five- or seven-tap SSEs, with the centre tap as the reference (zero index) tap to give equal effectiveness against minimum and non-minimum phase fading. The equalizer taps and weights are generally complex to enable the simultaneous equalization of the I and Q channels and interference between them (crossrail ISI).

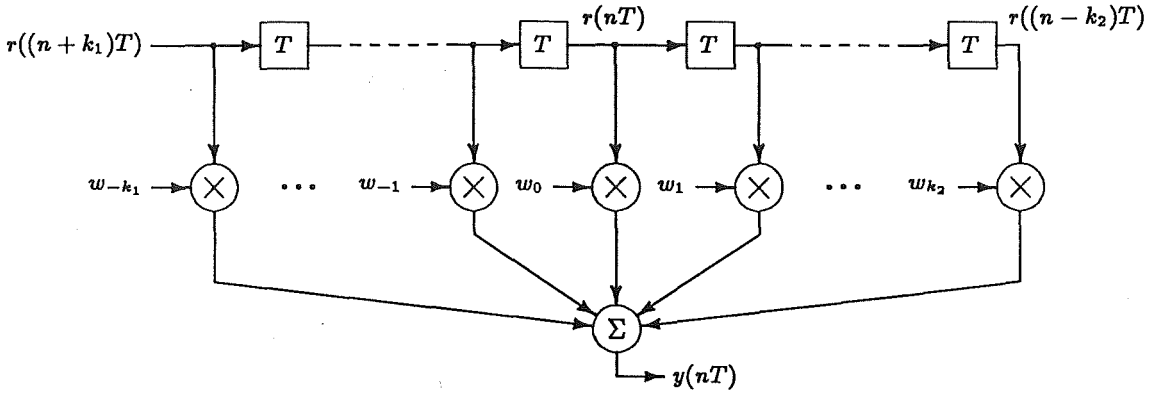


Figure 4.7 Synchronously-spaced equalizer.

Ideally the equalizer tap weights should be chosen to minimize the probability of a decision error in the receiver. This strategy, however, is difficult to implement in practice because the error probability is a highly nonlinear function of the tap weights. The more tractable criteria of minimum peak distortion or minimum mean-square error are usually used.

Lucky [1965] developed the *zero-forcing* (ZF) algorithm that computes tap weights to force the samples of the combined channel and equalizer impulse response to zero at all but one of the $(k_1 + k_2 + 1) T$ -spaced instants in the span of the equalizer. This minimizes the *peak distortion*

$$D_p \triangleq \frac{\sum_{i=-\infty}^{\infty} |\operatorname{Re}[h(iT)]| + \sum_{i=-\infty}^{\infty} |\operatorname{Im}[h(iT)]|}{|\operatorname{Re}[h(0)]|} \quad (4.10)$$

of the equalized impulse response $h(t)$, subject to the constraints of the equalizer length and delay. The ' on the summation indicates that the $i = 0$ term is not included.

The optimum tap weights for an SSE to equalize a response $h(t)$ using the ZF criterion are given by

$$\begin{bmatrix} h(0) & \cdots & h(-(k_1 + k_2)T) \\ h(T) & \cdots & h(-(k_1 + k_2 - 1)T) \\ \vdots & & \vdots \\ h(k_1 T) & \cdots & h(-k_2 T) \\ \vdots & & \vdots \\ h((k_1 + k_2 - 1)T) & \cdots & h(T) \\ h((k_1 + k_2)T) & \cdots & h(0) \end{bmatrix} \begin{bmatrix} w_{-k_1} \\ \vdots \\ w_{-1} \\ w_0 \\ w_1 \\ \vdots \\ w_{k_2} \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (4.11)$$

An infinite-length ZF equalizer is an inverse filter that inverts the folded frequency response of the channel. In practice, finite length equalizers that approximate the inverse

filter must be used. These may cause excessive *noise enhancement* at frequencies where the folded channel spectrum has high attenuation. The ZF criterion is only guaranteed to minimize peak distortion if the *peak ISI* before equalization is less than unity (eye open) [Lucky, 1965].

The *minimum mean-square error* (MMSE) criterion [Haykin, 1986] is more robust than the ZF criterion because it accounts for the noise by choosing tap weights to minimize the sum of squares of all ISI terms plus noise (mean-square error) at the output of the equalizer. Thus the *signal-to-distortion ratio* is maximized at the equalizer output, subject to the constraints of equalizer length and delay. The MMSE is a quadratic function of the tap weights; thus minimization of the mean-square error is always guaranteed, independent of noise or interference levels.

The optimum tap weights for an SSE to equalize a response $h(t)$ using the MMSE criterion are given by

$$\begin{bmatrix} r_{hh}(0) & \cdots & r_{hh}(-(k_1 + k_2)T) \\ r_{hh}(T) & \cdots & r_{hh}(-(k_1 + k_2 - 1)T) \\ \vdots & & \vdots \\ r_{hh}(k_1T) & \cdots & r_{hh}(-k_2T) \\ \vdots & & \vdots \\ r_{hh}((k_1 + k_2 - 1)T) & \cdots & r_{hh}(T) \\ r_{hh}((k_1 + k_2)T) & \cdots & r_{hh}(0) \end{bmatrix} \begin{bmatrix} w_{-k_1} \\ \vdots \\ w_{-1} \\ w_0 \\ w_1 \\ \vdots \\ w_{k_2} \end{bmatrix} = \begin{bmatrix} h^*(k_1T) \\ \vdots \\ h^*(T) \\ h^*(0) \\ h^*(-T) \\ \vdots \\ h^*(-k_2T) \end{bmatrix} \quad (4.12)$$

where $r_{hh}(jT) = \sum_i h(iT)h^*((i - j)T) + \sigma_\eta^2 \delta_{ij}$ and δ_{ij} is the Kronecker delta function

$$\delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases} \quad (4.13)$$

Noise enhancement can be further reduced with a *decision-feedback* equalizer (DFE) [Belfiore and Park, 1979] as shown in Figure 4.8. The forward filter tap weights reduce the contribution of the precursors to the ISI (using the MMSE or ZF criterion), and the feedback filter tap weights cancel the postcursors within the span of the filter. This equalizer is nonlinear because of the hard-decision element in the feedback path. Since the feedback filter does not enhance noise, the DFE can compensate for amplitude distortion with less noise enhancement than a linear equalizer. The feedback mechanism can, however, result in the propagation of decision errors, but this is not catastrophic and does not significantly degrade performance at usable error rates. Typical DFE's in DMR have two or three taps in the forward filter and one or two taps in the feedback filter.

DFEs perform well with minimum phase fading, when the post-cursor ISI is dominant, but their performance is similar to that of a linear equalizer with non-minimum phase fading, when the precursor ISI is dominant [Chamberlain *et al.*, 1986]. The overall performance of a DFE depends on the relative frequency of occurrence of minimum and non-minimum phase fading.

The frequency response of an equalizer with synchronous tap spacings is periodic with period $1/T$. Since Nyquist filters usually have 30 to 50% excess bandwidth, the sampled received signal is aliased. The effect of the aliasing is strongly dependent on the sampler phase [Qureshi, 1985]. In contrast, a *fractionally-spaced* equalizer (FSE) can offer improvements over a T -spaced equalizer. The FSE illustrated in Figure 4.9 has taps spaced at a fraction of a symbol period $T_s = KT/M$ where M and K are integers and $M > K$.

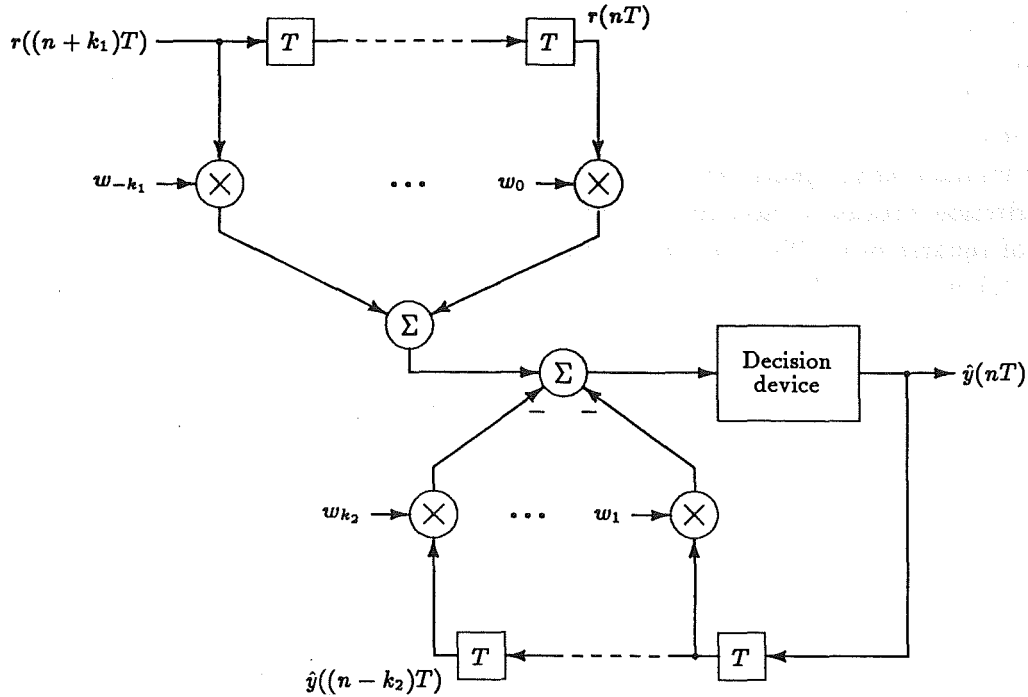


Figure 4.8 Decision-feedback equalizer.

Practically, it is convenient to choose $T_s = T/M$ where M is a small integer. An FSE can synthesize the best combination of adaptive matched filter and T -spaced equalizer within the constraints of its length and delay. The performance of an FSE is insensitive to sampler phase because the sampled received signal is no longer aliased (provided T_s is sufficiently small). It can also compensate for more severe delay distortion and deal with amplitude distortion with less noise enhancement than T -spaced equalizers. In practice, a $T/2$ -spaced equalizer performs as well as or better than a T -spaced equalizer with the same number of taps, even though the span of the FSE is half the span of the SSE. The MMSE criterion is usually used with an FSE.

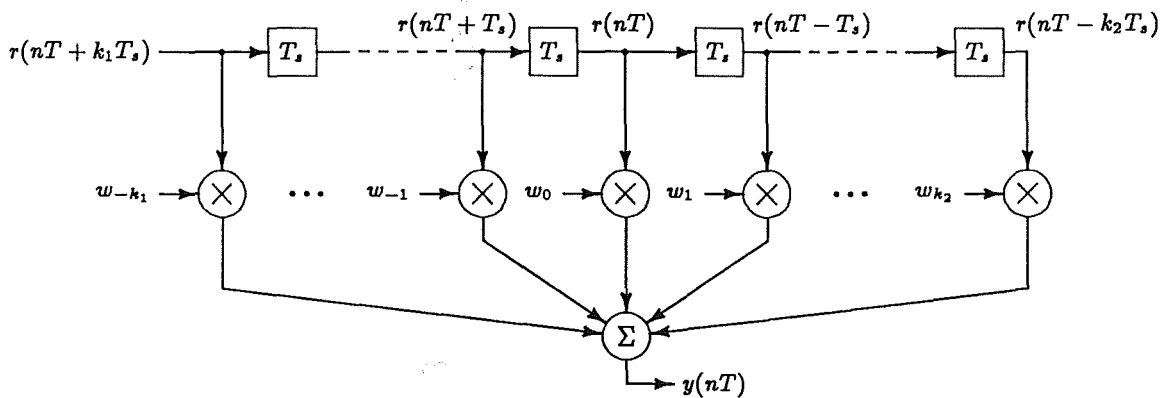


Figure 4.9 Fractionally-spaced equalizer.

The operation of a time-domain adaptive equalizer typically involves two modes. First the equalizer must be trained. Training involves transmitting a known sequence of symbols, which is matched against a synchronized version in the receiver to acquire information about the channel characteristics. The training sequence is usually a *pseudo-noise sequence*. The error between the training sequence and the equalizer output sequence is

used to adjust the tap weights in accordance with the ZF or MMSE criteria. After training, the equalizer is switched to a decision-directed mode, in which the difference between hard and soft decisions is used to adjust the tap weights. Adaptive equalizers are often used to compensate for modem imperfections such as non-ideal filters, modulators, and amplifiers in addition to fading.

Because the DMR channel only experiences time-dispersive fading for short periods of time, a formal training mode is not required. The system can be started under nominal propagation conditions, with the reference tap weight set to unity and all other tap weights set to zero.

Transversal equalizers can be realized before demodulation at IF, using an analog implementation, or after demodulation at baseband, using a digital implementation. The ZF algorithm is almost universally used for DMR applications because of its simple implementation.

Qureshi [1985] provides an excellent review of time-domain adaptive equalization techniques. The review is oriented to voice-band modem applications, but much of the material is also applicable to DMR systems.

4.3.2 Diversity

The essence of diversity is to exploit the dynamic, frequency selective, or spatially selective properties of multipath fading. Space diversity is the most common diversity technique for DMR systems, and takes advantage of the spatially selective nature of multipath fading. Multiple receive antennas are used with sufficient vertical separation so that the paths from the transmit antenna to each receive antenna are essentially independent. The receiver then selects the best signal out of the multiple received signals or combines the signals so that the transmission is more robust to the fading.

Frequency diversity takes advantage of the frequency selective nature of the fading. The same signal is transmitted on multiple carrier frequencies, and the receiver selects the best signal or combines the signals. Frequency diversity does not require additional antenna hardware, but does use bandwidth inefficiently. Lin *et al.* [1988] provide a review of diversity techniques for DMR systems.

4.3.3 Error-Control Coding

Error-control coding is particularly suitable for combatting errors due to the residual ISI after other countermeasures have been applied. It also provides resistance to errors induced by flat fading and lowers the background error rate under normal propagation conditions, when the channel is essentially an AWGN channel.

Conventional error correcting codes with code rate k/n require n/k times bandwidth expansion to accommodate the code redundancy. Consequently, low rate block and convolutional codes have not been widely used because they compromise bandwidth efficiency. High rate block and convolutional codes can provide about 3 dB coding gain without excessive erosion of spectral efficiency. Trellis-coded modulation (TCM), however, is particularly suited to bandwidth efficient applications because the information rate and bandwidth are not sacrificed to accommodate the code redundancy. TCM with Viterbi decoding can achieve coding gains of 3–6 dB on an AWGN channel and also achieve significant coding gains on dispersive channels (see Chapter 5). The implementation of a Viterbi decoder at DMR data rates is difficult with current technology, but will be possible in the near future using very large scale integration (VLSI) and parallel computing architectures.

QUALCOMM currently market a VLSI circuit containing a convolutional encoder and a Viterbi decoder capable of operating at 25 Mbits/s [QUALCOMM, 1990].

Error correcting codes can actually achieve larger performance gains than diversity techniques on channels experiencing normal propagation or flat fading. A diversity system with two diversity branches and *maximal-ratio combining* can provide a maximum of 3 dB gain on an AWGN channel. An error correcting code can provide 3–6 dB gain on the same channel, with a potentially lower cost.

4.4 Measures of Performance

The most useful measures of performance for digital communication systems are usually based on the probability of an error $P[\mathcal{E}]$ in a specified unit of data. The performance of a DMR system with normal propagation conditions or with flat fading depends only on the SNR of the received signal (neglecting modem imperfections). Curves of $P[\mathcal{E}]$ as a function of SNR are often used so that the distance between repeaters can be chosen to give an SNR that yields a desired $P[\mathcal{E}]$.

To measure the average performance of a system over a period of time, we must consider the duration of fading events and the statistics of the fade parameters. A threshold on bit error rate (BER) is set, for example $\text{BER} = 10^{-3}$, and when the system performance is worse than this threshold, the system is deemed to be in *outage*. The *outage time* is the accumulation of all periods of outage within a given time span—one year for example.

The International Telephone and Telegraph Consultative Committee (CCITT) has defined performance objectives in terms of budgets for *severely errored seconds*, *degraded minutes* and *error-free seconds* [CCITT, 1988]. In addition, a residual BER value must be respected. The severely errored seconds measure is defined as the percentage of seconds during which the BER exceeds 10^{-3} , and degraded minutes is defined as the percentage of minutes during which the BER exceeds 10^{-6} . Error-free seconds is defined as the percentage of whole seconds that are error free in a specified measurement interval covering many seconds. The residual BER is the BER measured in the absence of fading.

The most accurate way to estimate the outage of a DMR system on a particular path is to install the system on the path and measure the outage time. This method, however, is very expensive and time consuming because we must wait for fading to occur naturally. A more satisfactory technique for designing systems is to use a statistical model for the multipath channel. The channel can be simulated in hardware and the actual modem used, or the channel and modem can be simulated in computer software. To compute the system outage, the following procedure is used [Greenstein and Shafi, 1987]:

1. Choose a performance measure and threshold.
2. Choose a statistical channel model.
3. Find the region of channel model parameters Ω over which the performance is worse than the specified threshold.
4. Find the fraction of fade responses for which model parameters lie within Ω .

Step 4 can be performed by integrating the joint pdf of the model parameters over Ω or by generating many sets of model parameters using Monte Carlo techniques and counting the fraction of fades within Ω . This fraction is known as the *conditional outage probability*. The *outage time* in a year can be computed by multiplying the conditional probability of outage by the total number of fading seconds T_f per year. To measure the total number of

fading seconds per year, criteria must be established that reliably indicate the occurrence of fading. Usually, fluctuations in the power spectral density across the system bandwidth are monitored.

Lundgren and Rummler [1979] found that outage time was dominated by frequency selective fading rather than flat fading; consequently, they used a two-path fading model. With this approach, the outage region is defined by the system *signature curve* or *M-curve*, as shown in Figure 4.10. The AGC combats flat fading, so the two-path model is satisfactory in many cases. If flat fading is important, the outage region must be defined by signature surfaces.

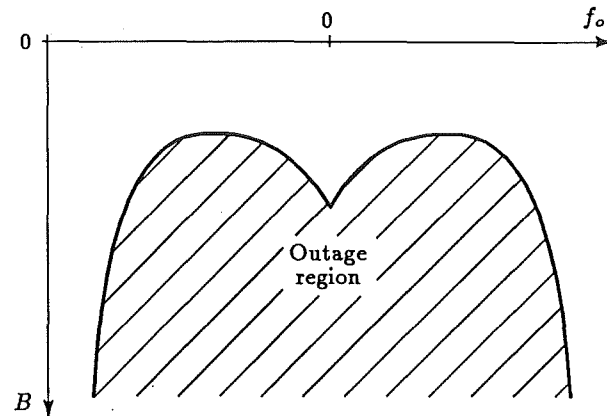


Figure 4.10 A typical system signature curve.

4.5 Conclusion

This chapter has discussed some important aspects of DMR systems. We will use the two-path Rummler model to describe the DMR channel. The DMR systems we will study use large signal constellations and incorporate equalization and TCM, but not diversity. The ISI that exists after non-ideal equalization is called *residual ISI*. We are interested in the performance improvements that TCM can offer on a channel that is experiencing residual ISI.

Chapter 5

Performance Estimation by Simulation

Simulation is a useful tool in the estimation of the performance of a digital communication system when a mathematical analysis is intractable, such as for coded systems on channels that introduce intersymbol interference (ISI). Monte Carlo simulation [Jeruchim, 1984] involves building a model of the system, usually in computer software, and estimating the probability of an event by empirical techniques. The model of the system must incorporate accurate statistical models of the sources of data and disturbances.

In this chapter we use Monte Carlo simulation to estimate the error rate for a unit of data at the receiver, when a block of data is transmitted. The size of the block of data governs the accuracy of the estimate. If the errors are mutually independent, the simulation time for a given accuracy is inversely proportional to the error rate. Hence, there is a lower limit to the error rates that can be estimated by simulation.

To measure low error rates ($< 10^{-2}$), Monte Carlo simulations must generate a large number of events that are not error events. Simulation times can be reduced, or conversely the lower limit on error rate can be reduced, by using *variance reduction* techniques [Jeruchim, 1984]. *Importance sampling* is a variance reduction technique that can be used in Monte Carlo simulations; it involves biasing the channel statistics so that the important events (those that cause errors) occur more often, and then unbiasing the estimated error rate to get the true error rate. Unfortunately, this technique is difficult to apply with Viterbi decoding because of the memory in the decoder. Herro and Nowack [1988] have used importance sampling to simulate a convolutional encoder with Viterbi decoding on an additive white Gaussian noise (AWGN) channel. Their simulation times are up to 20 times less than those for Monte Carlo simulation. Chen *et al.* [1990] have studied several variance reduction techniques for systems with memory. Some of these techniques show improvements over importance sampling.

We study the performance of digital microwave radio (DMR) systems with trellis-coded modulation (TCM), using Monte Carlo computer simulations. Because we can only use Monte Carlo simulation to estimate error rates down to about 10^{-5} , we will adopt an analytical approach to estimate lower error rates. The intention of the work reported in this chapter is to support the study of analytical performance bounds with ISI in Chapter 7.

A number of recent studies, independent of this study, have used simulation to examine the performance of TCM on channels that introduce ISI. Wong and McLane [1988] studied the application of TCM to high frequency radio channels. Their results show that TCM can achieve some performance improvement with the residual ISI introduced by this channel

after non-ideal equalization. Chouly and Sari [1988] examined the application of TCM to DMR systems without equalization. They show that coding gains are possible with the ISI introduced by light multipath fading, but the advantages of TCM are rapidly eroded with the deep fading often experienced by real systems. Therefore, TCM alone cannot be regarded as a countermeasure to fading, and any commercial DMR systems employing TCM must use adaptive equalizers, or other countermeasures, to significantly reduce the effects of frequency selective fading. Despinic *et al.* [1989] show that modest improvements in link outage can be obtained by applying TCM to a 256-QAM system with equalization.

We consider a 256-QAM DMR system with a data rate of 140 Mbits/s, and examine its performance when the data is trellis coded and mapped onto either a 512-CR or a 1024-QAM signal set [Carlisle *et al.*, 1990a]. The effects of different equalizers, codes, and decoding depths on performance are examined for a variety of multipath channels. Further specifications for the systems simulated are given in Section 5.1. The performance of the systems is measured by studying symbol error rate (SER) as a function of signal-to-noise ratio (SNR), and by estimating link outage. We are interested in the improvements that TCM can offer in terms of coding gain and reduction of the symbol error rate on a residual ISI channel, and how these improvements compare to those possible on an AWGN channel. These aspects of performance are studied in Section 5.2. We are also interested in the improvements that TCM can offer in terms of link outage on a residual ISI channel. The performance of DMR systems in terms of link outage is studied in Section 5.3.

5.1 Simulation Specifications

A general block diagram for the DMR systems that were simulated is given in Figure 5.1. Both coded and uncoded DMR systems were simulated for a variety of fading conditions. For the uncoded systems, the convolutional encoder and Viterbi decoder were removed, and the signal mapper was modified to perform Gray code mapping. Transmit and receive filters were root Nyquist filters with a combined frequency response from the raised cosine family. The roll-off factor of the raised cosine response was $\alpha = 0.3$. A Nyquist bandwidth of 17.5 MHz was used, this corresponds to an information rate of 140 Mbits/s for the 8 bits per symbol transmitted by the DMR systems.

The effect of multipath fading on the DMR channel was modelled using the Rummler model (see Section 4.2.1), with a delay between paths of $\tau = 6.3$ ns. Automatic gain control (AGC) was applied to restore signals to their unfaded power levels.

Three types of equalizer were used in the simulated systems—a synchronously-spaced equalizer (SSE), a fractionally-spaced equalizer (FSE) with a $T/2$ tap spacing, and a decision-feedback equalizer (DFE). All equalizers had five taps because this is a popular choice in commercial DMR systems. For the synchronously- and fractionally-spaced equalizers, the centre tap was the reference, and tap weights were optimized using the minimum mean-square error (MMSE) criterion. The DFE consisted of a forward filter with three taps (two precursor taps and a reference tap) and a feedback filter with two taps. The tap weights of the forward filter were optimized using either the MMSE or the zero-forcing (ZF) criterion, and the tap weights of the feedback filter were set to cancel the postcursor ISI within the span of the filter. The receiver sampling instants were adjusted to the maximal eye opening point using square law envelope detection prior to the equalizer, followed by narrowband filtering of the spectral line at the baud rate [Amitay and Greenstein, 1984].

An uncoded 256-QAM system was used as a reference system against which the coded

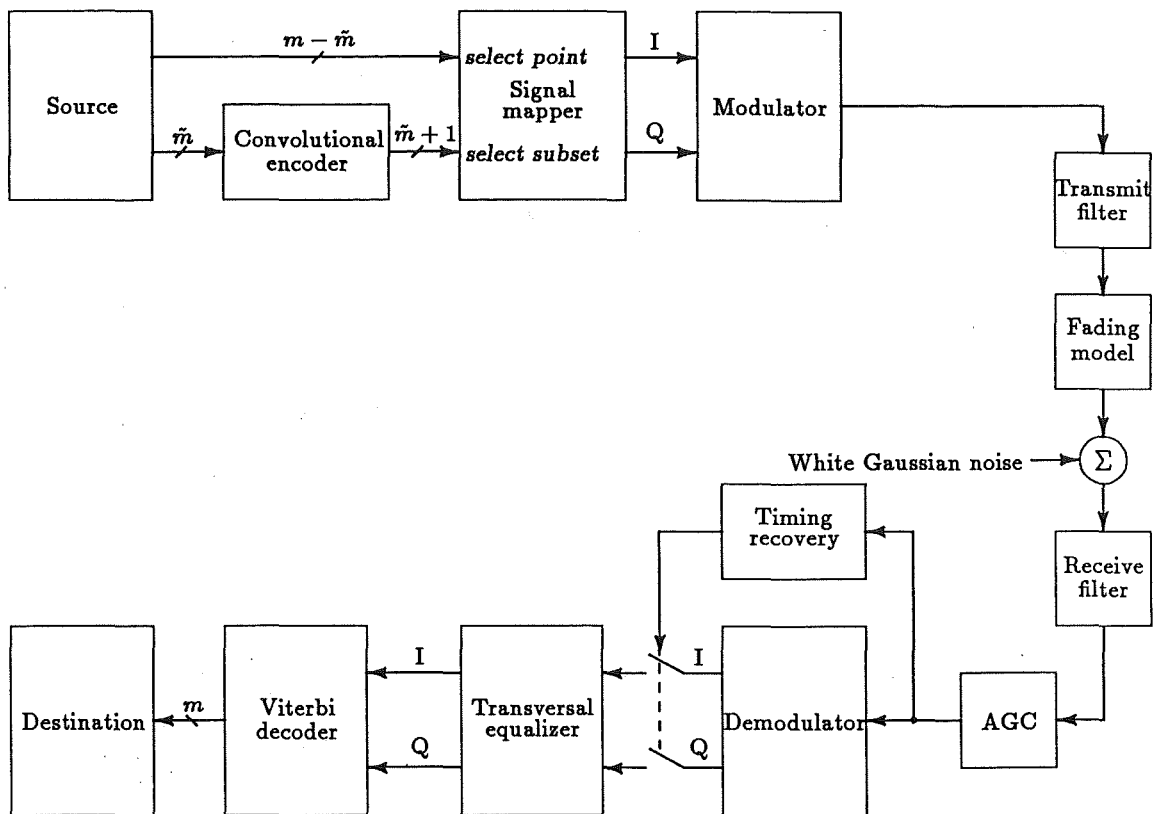


Figure 5.1 Block diagram of DMR system incorporating equalization and TCM.

systems could be compared. The bandwidth efficiency of this system was the same as the coded systems.

Two types of TCM scheme were studied. In one scheme, the in-phase (I) and quadrature (Q) channels were encoded separately as suggested by Thapar [1984]. Due to the reduced number of signal points per subset, compared to a standard Ungerboeck code, fewer metric calculations were required for subset decoding. This resulted in a significant reduction of computational cost compared to treating the channels together. A similar idea has recently been suggested by Viterbi *et al.* [1989] as a pragmatic approach to the implementation of TCM. The rate $1/2$, $\nu = 2$ (two binary memory elements) convolutional encoder shown in Figure 5.2a was used to add one redundant bit per symbol. Two of these encoders were used; one on the I channel and one on the Q channel. This resulted in a 1024-QAM coded signal constellation, which will be referred to as the *coded 1024-QAM system*. Separate Viterbi decoders were used to decode the I and Q channels.

The other TCM schemes used a single encoder for the I and Q channels, and mapped onto a 512-QAM constellation in the shape of a cross (512-CR). The rate $2/3$, $\nu = 3$ convolutional encoder shown in Figure 5.2b was used to add one redundant bit per symbol. This particular rate $2/3$ convolutional encoder is nonlinear and produces a TCM scheme that is invariant to carrier phase errors of 90° [Wei, 1984b]. To examine the effect of convolutional encoder memory on performance, the 512-CR signal constellation was also simulated with the rate $2/3$, $\nu = 5$ convolutional encoder shown in Figure 5.2c. These systems will be referred to as *coded 512-CR systems*.

The relative performance of the TCM schemes at high SNRs on an AWGN channel can be measured by comparing their asymptotic coding gains. Table 5.1 contains values

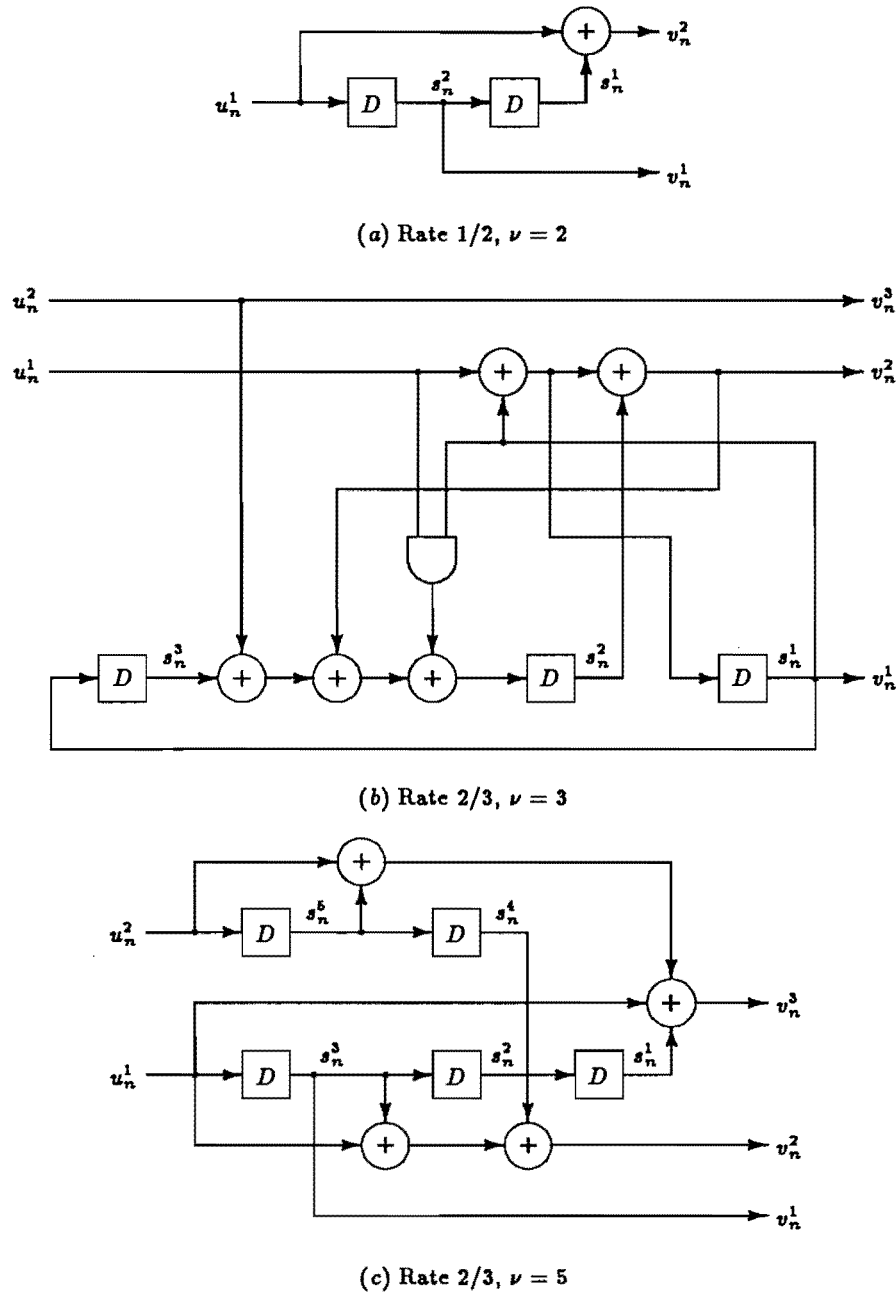


Figure 5.2 Convolutional encoders.

of asymptotic coding gain, computed using (3.7), for the various TCM schemes. Values of squared free distance d_{free}^2 and average energy E_s in the signal constellation are also given, where d is the minimum distance between signal points in the QAM constellations. Notice that $\nu = 5$ coded 512-CR has the highest asymptotic coding gain, and $\nu = 2$ coded 1024-QAM has the lowest asymptotic coding gain. Thus, on an AWGN channel, except at low SNRs, $\nu = 5$ coded 512-CR achieves the best performance and $\nu = 2$ coded 1024-QAM has the worst performance. We will see that this relative performance of the codes also holds with residual ISI.

The Viterbi decoder used a soft Euclidean metric (within the precision of the computer) on which to base its decisions. Forced decisions (see Section 3.3.1) were made by selecting the path with the lowest metric. In a real implementation, the Viterbi decoder would use quantization of eight levels between signal points, and forced decisions would

TCM	E	d_{free}^2	coding gain (dB)
uncoded 256-QAM	$42.5d^2$	d^2	—
$\nu = 2$ coded 1024-QAM	$170.5d^2$	$9d^2$	3.51
$\nu = 3$ coded 512-CR	$82.5d^2$	$5d^2$	4.11
$\nu = 5$ coded 512-CR	$82.5d^2$	$6d^2$	4.90

Table 5.1 Average signal point energy, squared free distance, and asymptotic coding gain for the various TCM schemes.

probably be made by selecting any path. As a result of these differences, we can expect our simulations to show a marginally better performance than might be expected with a real implementation. However, the method used is consistent with the analytical approach in Chapter 7.

Decoding depths of four to six constraint lengths were used in the Viterbi decoder (see Section 3.3.1). A decoding depth of twelve symbols was used for the coded 1024-QAM system. For the $\nu = 3$ and $\nu = 5$ coded 512-CR systems, decoding depths of twelve and twenty were used respectively; simulations with double these decoding depths produced increments in coding gain of less than 0.1 dB.

All results were obtained by computer simulations, using the Monte Carlo technique on equivalent lowpass representations (see Section 2.2) of signals and transfer functions. Thus the modulator and demodulator were not explicitly simulated. Symbol error rates were computed down to about 10^{-5} , which was the practical limit with the computing resources available. For an uncoded system with no ISI, the symbol errors are independent and binomially distributed. In this case, the variance of the error rate estimate is given by

$$\sigma^2(\hat{p}_e) = \frac{p_e(1 - p_e)}{N} \quad (5.1)$$

where p_e is the actual error rate, \hat{p}_e is the estimated error rate, and N is the number of trials. With the introduction of coding and ISI however, symbol errors are no longer independent, and the binomial distribution does not apply. The dependence between errors increases the variance of the error rate estimate. To obtain reasonable confidence on all results, within the constraints of computer time, the systems were simulated until at least one hundred symbol errors had been detected.

5.2 Error Rate Results

Simulations were performed to determine the coding gains and symbol error rates (SERs) of the various combinations of equalizer and TCM on an AWGN channel and frequency selective channels. To avoid simulating dynamic fading channels, three representative stationary fade characteristics conforming to the Rummmler model were considered, and the equalizer tap weights were held fixed at their optimum values throughout a simulation. The fading conditions selected were a spectral notch with $f_o = 0$ MHz (centred spectral notch), a spectral notch with $f_o = 4$ MHz (offset spectral notch), and a spectral slope across the system bandwidth. The depth of the centred spectral notch was selected to be 16 dB, which gave a residual distortion level (see (4.10)) of $D_p = 0.05$ following equalization by the five-tap SSE. The other fading conditions were selected to give the same residual

distortion with a five-tap SSE, resulting in an offset spectral notch of depth 9.4 dB and a spectral slope of 0.257 dB/MHz (4.5 dB amplitude variation over the Nyquist bandwidth).

The performance of the various DMR systems is calculated by adding white Gaussian noise to the faded channel and plotting SER against SNR. To specify the performance benefits of TCM quantitatively, we also calculate the coding gain for $\text{SER} = 10^{-4}$ and the reduction in SER for a given SNR.

In some of the results presented, the coding gain with residual ISI is greater than with AWGN alone. This cannot be interpreted as the code performing better with residual ISI than with AWGN alone. When the SER curves reach an error floor, it could be said that there is infinite coding gain, but the coded system is *not* performing infinitely better than the uncoded system. For these reasons, although coding gain is an appropriate measure of code effectiveness on an AWGN channel, it is better to consider gains in SER when the residual ISI is more severe (has greater variance) than the AWGN.

We define SER gain as the base ten logarithm of the ratio of uncoded SER to coded SER at a given SNR. A reference SER of 10^{-2} for the uncoded system has been used so that gains of up to two orders of magnitude (OM) can be measured from the results. It would be preferable to use a lower reference SER (e.g. 10^{-4}), but this would require SERs to be estimated down to about 10^{-7} , which is not possible with the resources available.

We consider each of the four channel conditions separately, and then discuss the overall trends observed. The simulation results are presented in the form of graphs of SER versus SNR, and specific performance measures are summarized in a table.

5.2.1 Additive White Gaussian Noise

The symbol error rate of the various TCM schemes on the AWGN channel is shown as a function of SNR in Figure 5.3. Table 5.2 summarizes the levels of SNR that yield $\text{SER} = 10^{-4}$ for the uncoded and coded systems, the coding gain (difference between the coded and uncoded SNRs), and the SER gain measured from Figure 5.3. The table also summarizes the asymptotic coding gains of the various TCM schemes relative to uncoded 256-QAM.

All the codes achieve coding gains at $\text{SER} = 10^{-4}$ that are in excess of 65% of the asymptotic coding gains. Despite having greater free distance (see Table 5.1), the coded 1024-QAM system achieves less coding gain and SER gain than either of the coded 512-CR systems. This is primarily due to the higher average energy in the 1024-QAM constellation. Incidentally, if the $\nu = 3$, rate $2/3$ convolutional encoder had been used in the coded 1024-QAM systems, the asymptotic coding gain would still have been 3.51 dB because the free distance is not improved. Thus, encoding the I and Q channels separately will reduce the complexity of the decoder, but the performance of the system is worse than that of similar

System TCM	SNR (dB) for $\text{SER} = 10^{-4}$			SNR (dB)	SER
	uncoded	coded	gain	asymptotic gain	gain
$\nu = 2$ coded 1024-QAM	31.3	29.1	2.2	3.51	0.8
$\nu = 3$ coded 512-CR	31.3	28.6	2.7	4.11	1.3
$\nu = 5$ coded 512-CR	31.3	28.1	3.2	4.90	2.0

Table 5.2 SNR levels, coding gains, and SER gains for the DMR systems on the AWGN channel.

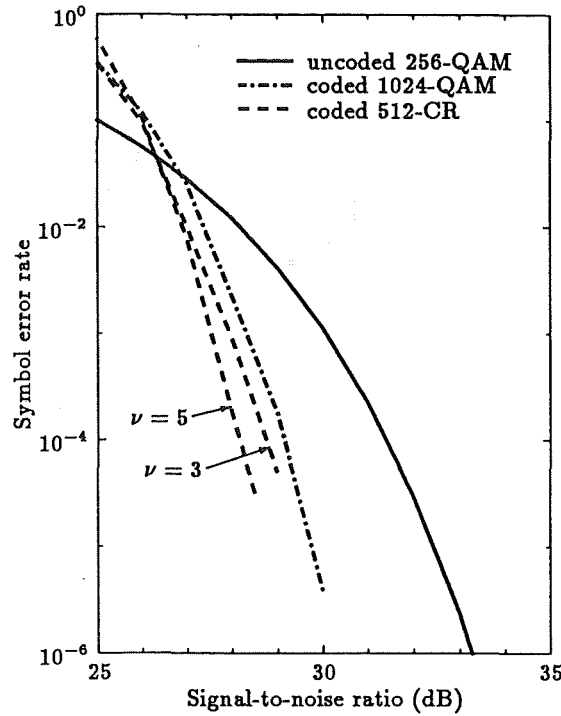


Figure 5.3 SER versus SNR curves for the DMR systems on the AWGN channel.

schemes that encode the I and Q channels together.

The effect of increasing the memory in the convolutional encoder for coded 512-CR from $\nu = 3$ to $\nu = 5$ is a 0.5 dB increase in coding gain and a 0.7 OM increase in SER gain.

5.2.2 Centred Spectral Notch

The symbol error rate of the various TCM schemes on the centred notch channel is shown as a function of SNR in Figures 5.4a and b. Table 5.3 summarizes the levels of SNR that yield $\text{SER} = 10^{-4}$ for the uncoded and coded systems, the coding gain, and the SER gain measured from Figure 5.4.

A comparison of SNR for the uncoded systems in Table 5.3 with SNR for the uncoded

System		SNR (dB) for $\text{SER} = 10^{-4}$			SER
equalizer	TCM	uncoded	coded	gain	gain
SSE	$\nu = 2$ coded 1024-QAM	35.5	32.4	3.1	0.7
	$\nu = 3$ coded 512-CR	35.5	31.7	3.8	1.3
	$\nu = 5$ coded 512-CR	35.5	31.0	4.5	2.0
FSE	$\nu = 3$ coded 512-CR	44.0	40.0	4.0	1.3
ZF-DFE	$\nu = 3$ coded 512-CR	37.0	32.8	4.2	1.0
MMSE-DFE	$\nu = 3$ coded 512-CR	—	—	$-\infty$	—

Table 5.3 SNR levels, coding gains, and SER gains for the DMR systems on the centred notch channel.

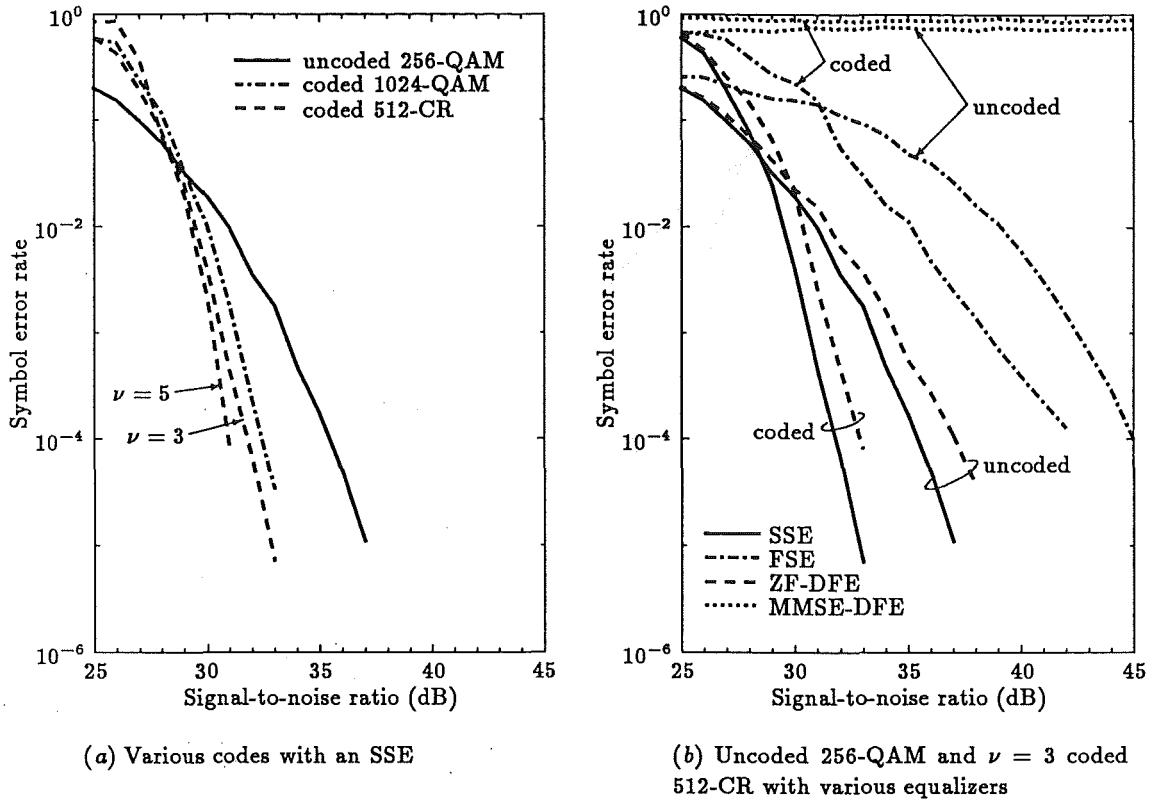


Figure 5.4 SER versus SNR curves for the DMR systems on the centred notch channel.

system on an AWGN channel gives a measure of the residual ISI and how well the equalizer is performing. Similarly, a comparison of SNR for a coded system with the same coded system on an AWGN channel gives a measure of how well the combination of equalizer and code is performing.

The coding gains for the codes with an SSE are all higher than for the AWGN channel, but the SER gains are very similar to those for the AWGN channel. The systems with an FSE do not perform well due to the short span of the equalizer and the lack of crossrail interference for the equalizer to work on. However, the $\nu = 3$ coded 512-CR still provides an SER gain of 1.3 OM. The SER gain of $\nu = 3$ coded 512-CR with the ZF-DFE is 0.3 OM less than with the SSE. Performance measures could not be determined for the MMSE-DFE because of the severely degraded performance compared to the other equalizers. The poor performance of the MMSE-DFE is due to attenuation of the cursor in the equalized impulse response, as discussed in Appendix 5A.

The effect of increasing the memory in the convolutional encoder for coded 512-CR from $\nu = 3$ to $\nu = 5$ is a 0.7 dB increase in coding gain and a 0.7 OM increase in SER gain.

5.2.3 Offset Spectral Notch

The symbol error rate of the various TCM schemes on the offset notch channel is shown as a function of SNR in Figures 5.5a and b. Table 5.4 summarizes the levels of SNR that yield $\text{SER} = 10^{-4}$ for the uncoded and coded systems, the coding gain, and the SER gain measured from Figure 5.5.

The coding gains for the codes with an SSE are again higher than for the AWGN channel, and the SER gains are very similar to those for the AWGN channel. The $\nu = 3$

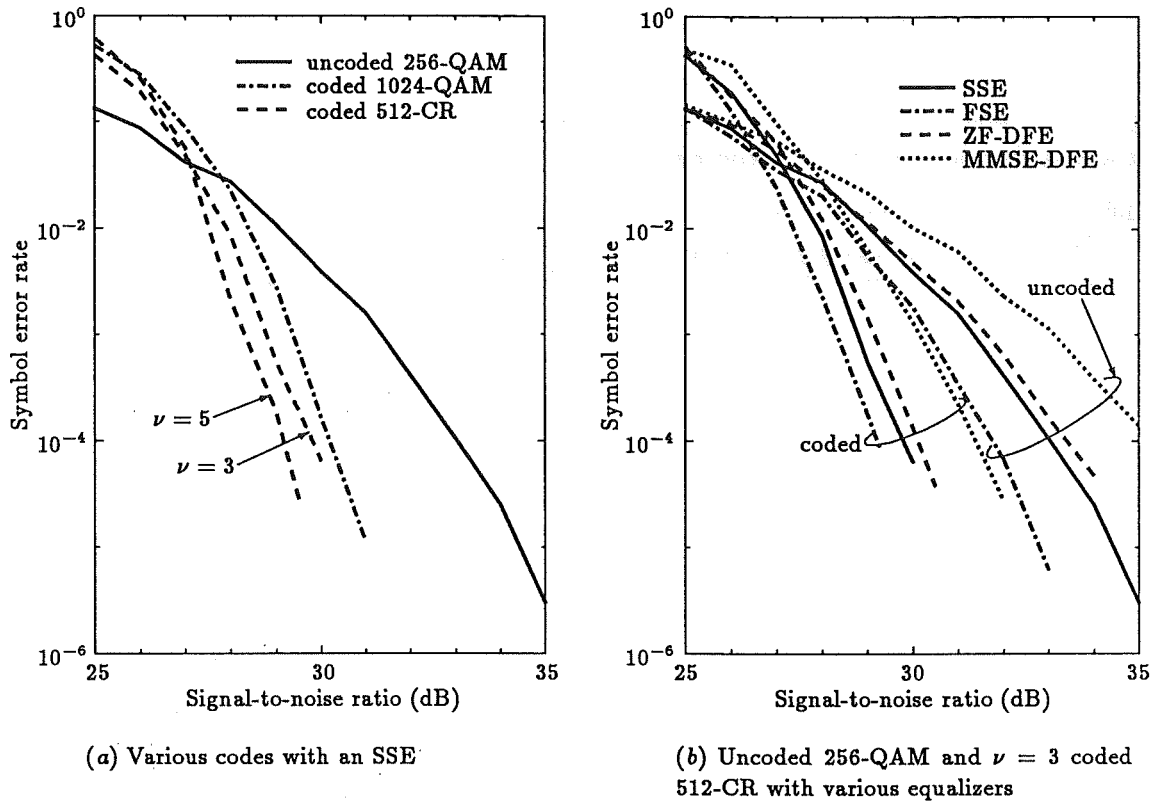


Figure 5.5 SER versus SNR curves for the DMR systems on the offset notch channel.

System		SNR (dB) for SER = 10^{-4}			SER gain
equalizer	TCM	uncoded	coded	gain	
SSE	$\nu = 2$ coded 1024-QAM	33.1	30.2	2.9	0.7
	$\nu = 3$ coded 512-CR	33.1	29.7	3.4	1.3
	$\nu = 5$ coded 512-CR	33.1	29.1	4.0	2.0
FSE	$\nu = 3$ coded 512-CR	31.8	29.2	2.6	1.2
ZF-DFE	$\nu = 3$ coded 512-CR	33.4	30.1	3.3	1.0
MMSE-DFE	$\nu = 3$ coded 512-CR	35.3	31.5	3.8	1.0

Table 5.4 SNR levels, coding gains, and SER gains for the DMR systems on the offset notch channel.

coded 512-CR in combination with the FSE achieves the best performance. The use of an MMSE-DFE again yields the worst performance due to cursor attenuation (Appendix 5A). The performance of the DFE is significantly improved by using the ZF criterion rather than the MMSE criterion. The SER gain of $\nu = 3$ coded 512-CR with either DFE is less than with the other equalizer types.

The effect of increasing the memory of the convolutional encoder for coded 512-CR from $\nu = 3$ to $\nu = 5$ is a 0.6 dB increase in coding gain and a 0.7 OM increase in SER gain.

5.2.4 Spectral Slope

The symbol error rate of the various TCM schemes on the slope channel is shown in Figures 5.6a and b. Table 5.5 summarizes the levels of SNR that yield $\text{SER} = 10^{-4}$ for the uncoded and coded systems, the coding gain, and the SER gain measured from Figure 5.6.

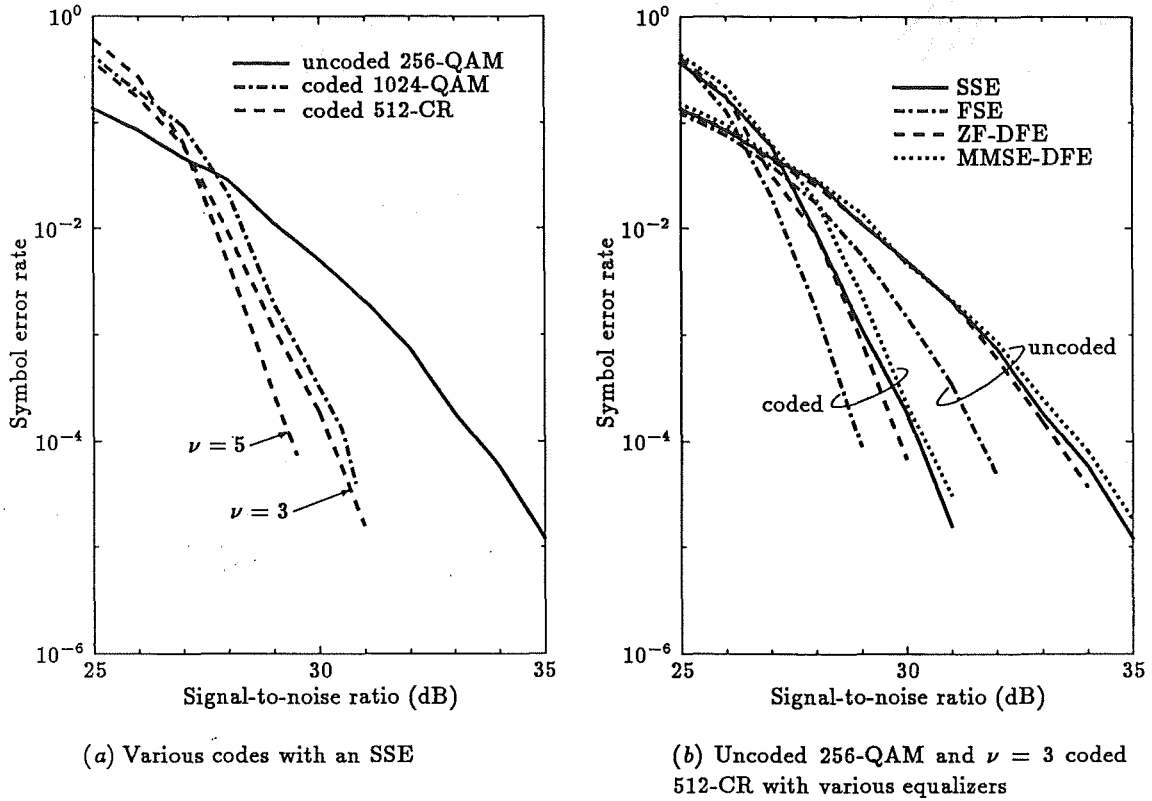


Figure 5.6 SER versus SNR curves for the DMR systems on the slope channel.

System		SNR (dB) for $\text{SER} = 10^{-4}$			SER
equalizer	TCM	uncoded	coded	gain	gain
SSE	$\nu = 2$ coded 1024-QAM	33.4	30.6	2.8	0.8
	$\nu = 3$ coded 512-CR	33.4	30.2	3.2	1.0
	$\nu = 5$ coded 512-CR	33.4	29.3	4.1	1.7
FSE	$\nu = 3$ coded 512-CR	31.7	29.0	2.7	1.4
ZF-DFE	$\nu = 3$ coded 512-CR	33.4	29.8	3.6	1.2
MMSE-DFE	$\nu = 3$ coded 512-CR	33.8	30.3	3.5	0.9

Table 5.5 SNR levels, coding gains, and SER gains for the DMR systems on the slope channel.

The trends observed in these results are very similar to those observed for the offset notch. The exception is that the SER gains of the coded 512-CR schemes with an SSE are degraded compared to the SER gains on an AWGN channel. This suggests that the coded 512-CR performs better with the residual ISI from a notch equalized with an SSE, than with the residual ISI from a slope equalized with an SSE.

The effect of increasing the memory of the convolutional encoder for coded 512-CR

from $\nu = 3$ to $\nu = 5$ is a 0.9 dB increase in coding gain and a 0.7 OM increase in SER gain.

5.2.5 Discussion of Results

For all the cases of residual ISI considered, the coding gain at the reference level of $\text{SER} = 10^{-4}$ was greater than or equal to the coding gain for an AWGN channel. In some cases, the coding gain with residual ISI was up to 55% larger than with AWGN. However, coding gain is a performance measure that is most appropriate for AWGN channels, and we cannot conclude that TCM offers greater performance improvements with residual ISI than with AWGN alone. In contrast with the coding gain results, the SER gains with residual ISI are generally very similar to those with AWGN. This suggests that TCM offers similar performance gains with residual ISI to those with AWGN alone.

The improvements in coding gain when the memory in the convolutional encoder is increased from $\nu = 3$ to $\nu = 5$ are primarily due to the larger free distance of the $\nu = 5$ code. But when the additional memory increases the constraint length of the encoder, it will also offer greater resistance to burst errors.

The SER gains are somewhat independent of the residual ISI characteristics, indicating that the codes are just as effective with the combination of residual ISI and AWGN as they are with AWGN alone. In general, the SER at the 10^{-2} reference level is improved by one to two orders of magnitude by the addition of coding. The SER gain results would be more useful if we could measure them at a lower SER, but this would require impractically long simulation times.

Although coding gain indicates the code performance for a given equalizer, it does not measure the overall performance of the code and equalizer in combination. The overall performance can be measured by comparing the SNRs of the various systems for $\text{SER} = 10^{-4}$. Except for the centred notch, we find that the FSE performs with the lowest SNR, even though it spans only 2.5 symbols with its five taps. In the case of the centred notch, the performance of the FSE degrades in comparison to longer-span synchronously-spaced equalizers. An FSE with a five-symbol span cannot perform worse than any of the other five-tap linear equalizers, and will generally perform significantly better than all other five-tap equalizers.

From the simulations in this section, we conclude that coding is able to reduce the error rate following imperfect equalization. To measure the average performance over a full range of fading conditions, we now compute link outage probabilities.

5.3 Outage Probability Results

The International Telephone and Telegraph Consultative Committee (CCITT) defines link outage in terms of budgets for severely errored seconds, degraded minutes, and error-free seconds [CCITT, 1988]. A residual BER value must also be respected. The definition of each of these terms is given in Section 4.4. Note that in the last section the performance measures were based on SER, but that in this section they are based on BER. The BER is always less than the SER; however, the exact relationship between BER and SER depends on the signal mapping, for an uncoded system, and the TCM scheme, for a coded system.

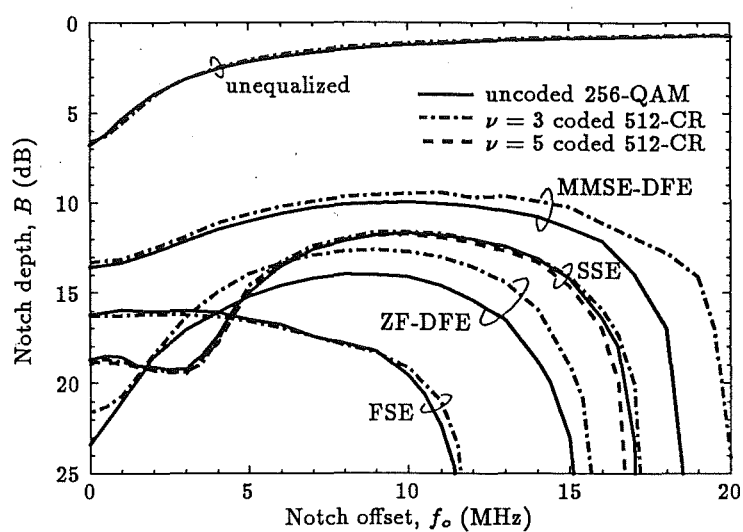
To investigate the improvements in outage possible with TCM, signatures for specific bit error rates of 10^{-3} and 10^{-4} were calculated by finding the fade depth that produced these error rates at notch frequencies (f_o) positioned inside and outside the system bandwidth. *Relative* measures of severely errored seconds and degraded minutes for systems

experiencing multipath fading can be estimated from the system signatures at $\text{BER} = 10^{-3}$ and $\text{BER} = 10^{-6}$ respectively. Unfortunately, $\text{BER} = 10^{-6}$ values require excessive computer time to simulate, so signatures were calculated for $\text{BER} = 10^{-3}$ and $\text{BER} = 10^{-4}$ thresholds. Trends in performance will be deduced from these results.

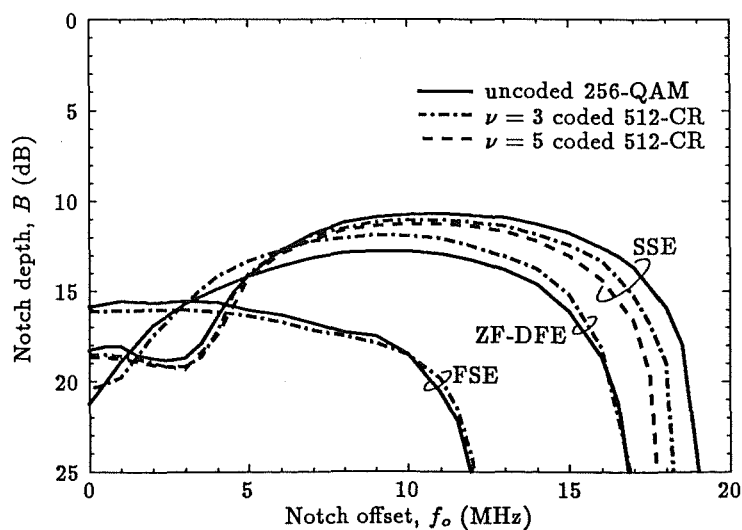
The error-free seconds and residual BER measures cannot be verified directly because they involve error rates much lower than those that can be estimated by simulation. However, from the the known characteristics of TCM it is possible to infer the effect of TCM on these measures.

5.3.1 System Signatures

All signatures were obtained by iteratively adjusting the notch depth at a given notch offset frequency until the threshold BER was obtained at the receiver. An SNR of 60 dB, which is a commonly used nominal value in the absence of flat fading, was maintained. The signatures for $\text{BER} = 10^{-3}$ are given in Figure 5.7a, and those for $\text{BER} = 10^{-4}$ are



(a) $\text{BER} = 10^{-3}$



(b) $\text{BER} = 10^{-4}$

Figure 5.7 System signatures.

given in Figure 5.7*b*. Since they are symmetrical about $f_o = 0$ MHz, only half of the signatures are shown.

As shown in Figure 5.7*a*, the unequalized signatures are poor for both the coded and uncoded systems, thus verifying that TCM alone is not a useful countermeasure to multipath fading [Chouly and Sari, 1988].

The contrasting performance of the different types of five-tap equalizers is made immediately apparent by the signatures. The SSE shows good tolerance to notches in the range $f_o = 0$ to 4 MHz, but the tolerance of the FSE to in-band slope is reflected in its superior performance for $f_o > 4$ MHz. Within the range $f_o = 0$ to 1 MHz, the ZF-DFE signatures are the best, but overall are only marginally better than those for the SSE. The signatures for the systems with an MMSE-DFE are worse than the signatures of all the other equalized systems for all notch offset frequencies.

Comparison of the signatures for a given equalizer with and without coding leads to the conclusion that, despite the significant coding gains found previously, the improvements in outage offered by TCM are relatively small at BER levels of 10^{-3} and 10^{-4} . A quantitative evaluation of this observation was made by computing outage probabilities from the system signatures.

5.3.2 Outage Computation from Signatures

The probability of outage, given that multipath fading is occurring, is known as the *conditional probability of outage*

$$P_o = \int \int_{\Omega} p(B)p(f_o) dB df_o \quad (5.2)$$

where Ω is the system signature for a given level of BER, and $p(B)$ and $p(f_o)$ are the probability density functions of the notch depth and frequency parameters during fading. In practice, multipath fading is deemed to be occurring when the power spectral density, measured across the system bandwidth at the receiver, fluctuates more than a predetermined level.

To obtain realistic estimates of outage, typical probability density functions $p(B)$ and $p(f_o)$ were formed from propagation measurements made during periods of multipath fading on a 60 km over-water path in New Zealand [McKay and Shafi, 1988]. These probability density functions are shown in Figures 5.8*a* and *b*. Numerical integration techniques were used to evaluate (5.2).

Results of the conditional outage probability computations for $\text{BER} = 10^{-3}$ are given in Table 5.6, which presents the conditional outage probabilities and the percentage improvements in outage due to coding. Comparison of the conditional outage probability figures for the uncoded systems shows that the outage of the systems with SSE, FSE and ZF-DFE equalizers will be improved by factors of 17, 33, and 26 over an unequalized channel. Thus, the equalizer is an important countermeasure for reducing outage. Again, although it only has half the symbol span of the other equalizers, the FSE performs the best.

Conditional outage probabilities and percentage improvements in outage for $\text{BER} = 10^{-4}$ are presented in Table 5.7. For the particular probability distributions of the fade parameters that have been assumed, and for $\nu = 3$, the percentage reduction in outage by incorporating coding is greatest for the FSE. Increasing the memory in the convolutional encoder to $\nu = 5$, in combination with an SSE, gives a marked improvement in outage. It is expected that this improvement will carry over to the $\nu = 5$ coded 512-CR with any of the equalizers, including the FSE, which will thus maintain its superiority.

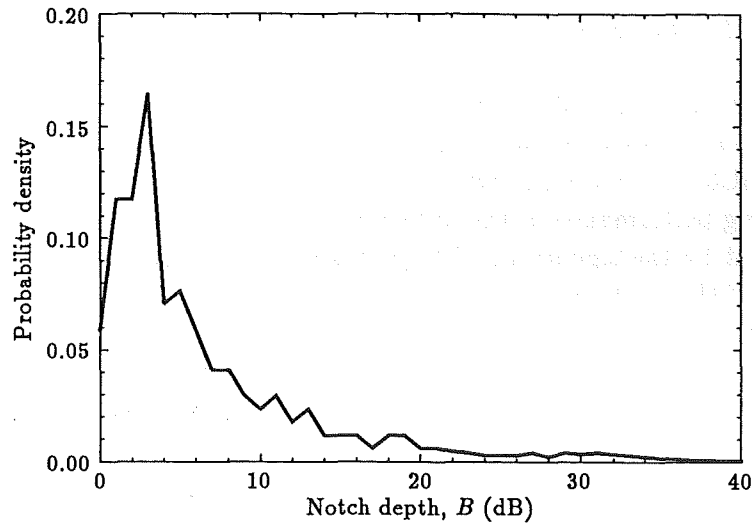
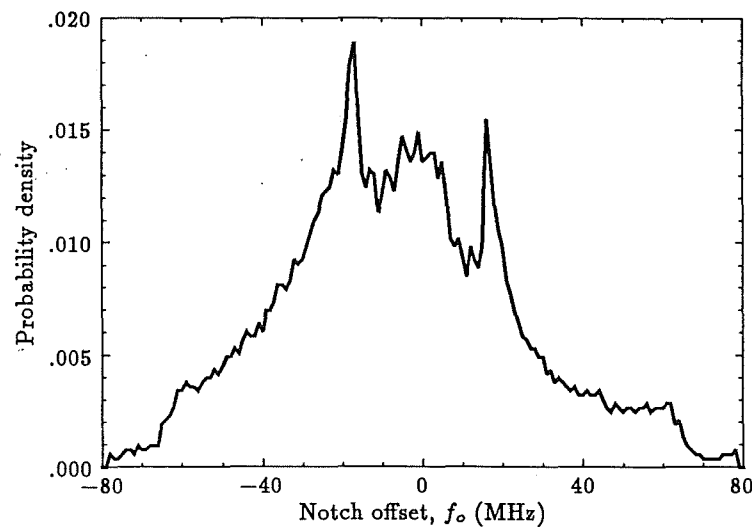
(a) Probability density function of notch depth B (b) Probability density function of notch offset f_o

Figure 5.8 Probability density functions of fading model parameters for an over-water path in New Zealand.

From the figures given in Table 5.6, the improvement in severely errored seconds (reference level of $\text{BER} = 10^{-3}$) brought about by coding will be insignificant. Although results could not be obtained for $\text{BER} = 10^{-6}$, the trends of improved outage due to coding as BER goes from 10^{-3} to 10^{-4} are expected to continue, giving useful gains in the degraded minute objective.

Despinic *et al.* [1989] have independently estimated outage improvements of up to 28% due to coding, when $\text{BER} = 10^{-3}$ for $\nu = 3$ coded 512-CR with a five-tap FSE. There are two possible explanations for this result being larger than the results presented here. First, they considered outage due to flat and dispersive fading, whereas only dispersive fading has been considered in this study. Second, their technique for computing conditional outage probabilities was different to the technique just described, and different fade parameter statistics were used.

System configuration		Conditional outage	Coding outage
Equalizer	Modulation	probability	improvement
none	uncoded 256-QAM	8.83×10^{-1}	reference
	$\nu = 3$ coded 512-CR	8.92×10^{-1}	-1.0%
SSE	uncoded 256-QAM	5.21×10^{-2}	reference
	$\nu = 3$ coded 512-CR	5.46×10^{-2}	-4.8%
	$\nu = 5$ coded 512-CR	4.89×10^{-2}	+6.1%
FSE	uncoded 256-QAM	2.48×10^{-2}	reference
	$\nu = 3$ coded 512-CR	2.44×10^{-2}	+1.6%
ZF-DFE	uncoded 256-QAM	3.42×10^{-2}	reference
	$\nu = 3$ coded 512-CR	4.52×10^{-2}	-32.2%
MMSE-DFE	$\nu = 3$ uncoded 256-QAM	8.53×10^{-2}	reference
	$\nu = 3$ coded 512-CR	9.46×10^{-2}	-10.9%

Table 5.6 Conditional outage probabilities and outage improvements due to coding for BER = 10^{-3} .

System configuration		Conditional outage	Coding outage
Equalizer	Modulation	probability	improvement
SSE	uncoded 256-QAM	6.95×10^{-2}	reference
	$\nu = 3$ coded 512-CR	6.45×10^{-2}	+7.2%
	$\nu = 5$ coded 512-CR	5.70×10^{-2}	+18.0%
FSE	uncoded 256-QAM	2.86×10^{-2}	reference
	$\nu = 3$ coded 512-CR	2.57×10^{-2}	+10.1%
ZF-DFE	uncoded 256-QAM	4.84×10^{-2}	reference
	$\nu = 3$ coded 512-CR	4.92×10^{-2}	-1.7%

Table 5.7 Conditional outage probabilities and outage improvements due to coding for BER = 10^{-4} .

5.3.3 Error-Free Seconds and Residual Bit Error Rate

The error-free seconds (EFS) objective is specified as 99.68% EFS for a 2500 km, 64 kbits/s reference channel [CCITT, 1988]. For a hop length of 50 km, and assuming a linear apportionment of the EFS objective, this translates to 99.9936% EFS for each hop. To estimate the corresponding symbol error rate, we assume a binomial distribution for symbol errors, so that the symbol error probability and percentage error-free seconds are related by

$$\text{EFS} = 100(1 - \text{SER})^{R_s} \quad (5.3)$$

where R_s is the symbol rate. This estimate of EFS is pessimistic because independent error events, rather than bursts, have been assumed. For 8 bit source symbols (256-QAM) on a 64 kbits/s channel, $R_s = 8000$ symbols per second, so that (5.3) gives $\text{SER} = 8 \times 10^{-9}$ at the target EFS. This corresponds to $\text{BER} \approx 10^{-9}$ when the symbols are Gray-encoded.

For $\text{BER} = 10^{-9}$, it is impractical to use simulation to study the gains in the EFS

margin due to the coding. Indeed, the measurements to verify that a real system is meeting the EFS objective require many months of monitoring. When $\text{BER} = 10^{-9}$, however, the coding gain is close to its asymptotic value, so the TCM will result in a significantly lower BER and help to meet the EFS objective.

This reduction in error rate, due to TCM, will also significantly increase the margins by which DMR systems meet the residual BER objective for a 2500 km reference channel.

5.4 Conclusion

Simulation has been used to study the improvements that TCM can make to the performance of an equalized 256-QAM digital microwave radio system. Although coding gain was achieved with $\nu = 2$ coded 1024-QAM on an AWGN channel, useful additional coding gain of 0.5 dB was obtained with $\nu = 3$ coded 512-CR. A further 0.5 dB of coding gain was achieved by increasing the memory of the code to $\nu = 5$.

The ability of TCM to compensate for residual ISI was demonstrated by the coding gain being maintained or increased, over that for AWGN alone, for a given SER in a residual ISI/AWGN environment. The SER gain at a reference level of 10^{-2} , for a given code, was approximately constant for the residual ISI and AWGN channels, indicating that TCM is just as effective on a residual ISI channel as it is on an AWGN channel. The results also indicate that TCM combined with fractionally-spaced equalization offers better overall performance than TCM with synchronously-spaced or decision-feedback equalizers. Optimum decoding depths for a residual ISI channel were found to be about the same as those for an AWGN channel.

Link outage is an important measure of DMR system performance. TCM with an FSE had the lowest outage of the systems simulated, at both the 10^{-3} and 10^{-4} BER thresholds, and showed the largest improvements for a given code. The use of TCM does not always result in an improvement to the severely errored seconds component of outage ($\text{BER} = 10^{-3}$); however, at $\text{BER} = 10^{-4}$ improvements in outage do occur. The trends lead us to speculate that these improvements will be even greater at $\text{BER} = 10^{-6}$ (degraded minutes), and that trellis-coded systems can achieve better error-free second and residual BER performance than equivalent uncoded systems. Codes with increased memory can be expected to provide further improvements in outage, but achievement of the improvements with coding will require first level countermeasures by adaptive equalization.

The results of this study are encouraging because they demonstrate that TCM can be used as a countermeasure for residual ISI on DMR systems, thereby improving the system performance. It should be noted that the effects of carrier recovery errors, timing jitter, high power amplifier nonlinearities, and roll-off factors, on the performance of TCM in a residual ISI environment, were not included in this study. Unfortunately, simulation is too slow to make a thorough study of all these factors without using many months of computer time. Also, the time constraint of simulation does not allow error rates below about 10^{-5} to be estimated. Tight analytical bounds on the performance of TCM for ISI channels are needed so that time consuming simulation can be avoided. This issue is addressed in Chapter 7.

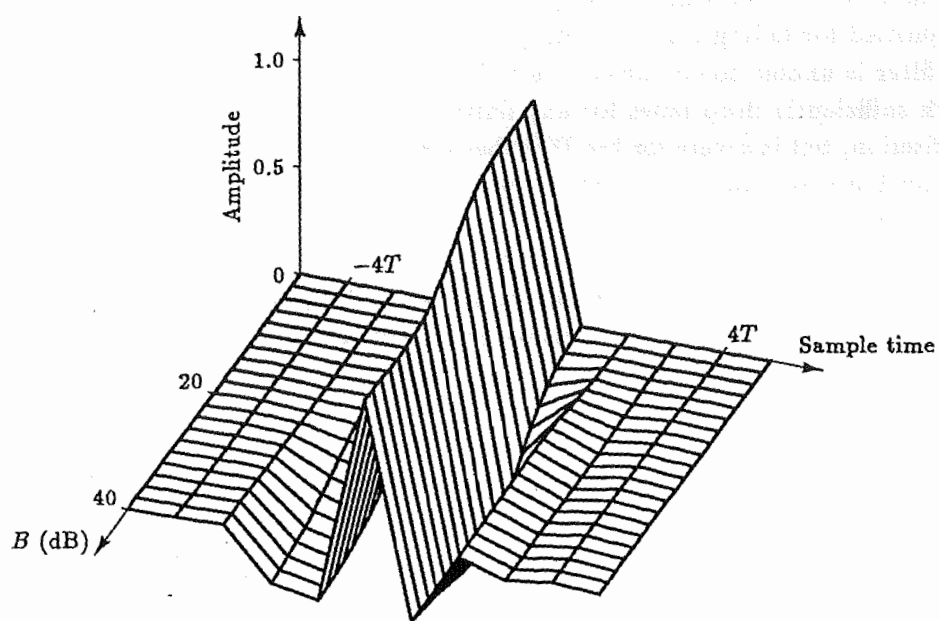
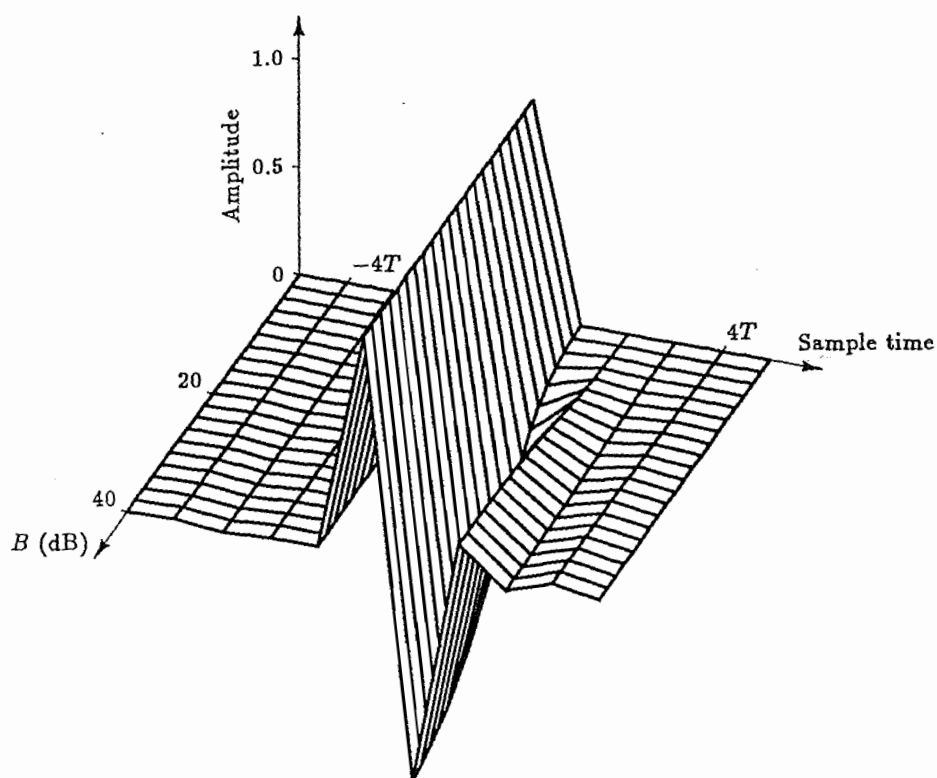
Appendix 5A Cursor Attenuation

The poor performance of the MMSE-DFE merits special attention. It is due to attenuation of the cursor in the equalized impulse response because the MMSE algorithm does not fully

remove the multiplicative interference. This phenomenon has been referred to as *cursor attenuation* [Carlisle *et al.*, 1989]. Cursor attenuation with the MMSE algorithm, at high SNR, is illustrated in Figure 5.9a. Sampled impulse responses at the output of the forward filter are plotted for fading consisting of centred spectral notches of different depths. The feedback filter is unable to compensate for the attenuation of the cursor. This effect will occur with sufficiently deep fades for any finite transversal equalizer optimized using the MMSE criterion, but is severe for the DFE because of the lack of taps in the forward filter to cancel postcursors. Also note that the precursor samples are just as significant as the postcursor samples.

Cursor attenuation causes errors when large signal levels are being sent. For example, the largest level in a 512-QAM constellation is $23d/2$, where d is the minimum distance between signal points, and a hard-decision error will result if this level is attenuated to less than $11d$. This occurs if the cursor in the equalized impulse response is attenuated to less than $22/23$, which is the case for centred notches deeper than 17 dB. Cursor attenuation is only a problem when the size of the signal constellation is large. Systems using 16-QAM constellations can tolerate fades of up to 40 dB without cursor attenuation being a major problem.

To avoid cursor attenuation, the ZF criterion can be used to optimize the forward tap weights in the DFE. Figure 5.9b shows the sampled impulse responses at the output of the forward filter, again with fading consisting of centre band notches of different depths. The ZF criterion forces the cursor in the impulse response at the output of the forward filter to unity, and concentrates the distortion in the first postcursor outside the span of the forward filter. The ISI due to the concentrated postcursor is then cancelled by the feedback filter, resulting in an improved overall response.

(a) MMSE-DFE ($S/N \rightarrow \infty$)

(b) ZF-DFE

Figure 5.9 Impulse response at the output of the forward filter of the DFE as a function of fade depth.

Chapter 6

The Probability Density of Intersymbol Interference

To compute decision error probabilities and to better understand the mechanisms causing errors, it is useful to know the *probability density function* (pdf) of the intersymbol interference (ISI) for a system on a specific channel [Hill, 1971; Metzger, 1987; Carlisle *et al.*, 1990b]. In this chapter we derive algorithms to compute approximate ISI pdf's for uncoded and trellis-coded systems on time-dispersive channels. These approximate ISI pdf's will be used in the next chapter to formulate analytical bounds on the error probability of TCM on time-dispersive channels.

For all but the simplest systems, on channels with impulse responses of short time duration, it is not usually feasible to compute the ISI pdf exactly. One exception appears in the work of Hill and Blanco [1973], where a tractable analytical expression for the ISI pdf is developed for uncoded binary transmission (they also note the extension to 2^m -ary PAM) on channels with infinite impulse responses expressible as geometric sequences of the form

$$h_i = \begin{cases} (1 - \beta)\beta^i & \text{if } i \text{ is a positive integer} \\ 0 & \text{otherwise} \end{cases} \quad (6.1)$$

where $\beta = 2^{-1/K}$ for any positive integer K . When $\beta \in (0, 1)$ the distribution function of z is either absolutely continuous or purely singular, so that the ISI pdf has no discrete components.

For channels with finite impulse responses, the ISI pdf consists of clusters of discrete components (mass points) and should, strictly speaking, be referred to as a *probability mass function* (pmf). In such cases researchers usually resort either to approximating ISI pdf's or to finding ISI pdf's that bound the symbol error probability. Huzii and Sugiyama [1970] also consider geometric impulse responses with finite state Markov pulse train transmission. The distribution function of the ISI is found by solving a set of functional equations by successive approximation. Hill [1971] takes a less accurate but more general approach and computes an approximate ISI pdf for 2^m -ary PAM constellations with partial response coding on a general channel. The ISI pdf is computed by convolving a number of partial pdf's. The computational cost is greatly reduced by quantizing to discrete levels of ISI and by decomposing the 2^m -ary constellation into a weighted sum of binary constellations. This technique can also be applied to 2^{2m} -ary QAM constellations with square boundaries by decomposing the constellation into a weighted sum of 4-QAM sub-constellations. Metzger [1987] uses different terminology but essentially uses

a variation on Hill's approach to approximate the ISI pdf.

Glave [1972] has computed distribution functions for the ISI that maximize the symbol error probability, subject to peak ISI and ISI-variance constraints. These distributions are only valid for sufficiently high signal-to-noise ratio (SNR). Matthews [1973] extended Glave's work to any SNR, and also computed distribution functions that minimize the symbol error probability subject to the same constraints.

In this chapter we are specifically concerned with computing ISI pdf's for high capacity *digital microwave radio* (DMR) systems employing equalization and *trellis-coded modulation* (TCM). These systems may have large, possibly non-square constellations (e.g. cross constellations) with decision regions that are not necessarily square (e.g. honeycomb constellations). Also, TCM introduces dependence between transmitted signal points. For these two reasons, the work of Hill [1971] and Metzger [1987] is not directly applicable.

Some preliminary notation is introduced in Section 6.1. Using this notation, algorithms are presented for computing approximate ISI pdf's for uncoded systems in Section 6.2. New algorithms are developed in Section 6.3 to evaluate the ISI pdf for systems that employ Ungerboeck codes, although the results are applicable to trellis-codes in general. In Section 6.4, ISI pdf's that provide definite lower and upper bounds on symbol error probability are examined. Finally, in Section 6.5 we study the effect on the ISI pdf's of parameters (e.g. level of quantization) used in the computation. We also examine the ISI pdf for coded signal constellations, and compare it to the ISI pdf for the same signal constellations without coding. The effect on the ISI of the dependence between signal points, introduced by a coded signal constellation, is thus studied. Note that since TCM cannot be expected to perform well with severe or *raw ISI* (i.e. without equalization), we will be primarily interested in determining the probability densities of the *residual ISI* after non-ideal equalization (i.e. with finite-tap equalizers). However, the theory and algorithms developed are equally applicable to all levels of ISI. The term ISI will be used to refer to both raw and residual ISI.

6.1 Preliminaries

ISI is manifest after sampling a received signal, so it is helpful to model the continuous-time channel by a discrete-time channel. The channel may be time varying, but will usually be stationary over a number of symbol intervals. In such cases the channel can be considered at a specific instant in time. We use the channel model shown in Figure 6.1 where, at time n , x_n is the transmitted signal point, z_n is the ISI, η_n is a sample from a Gaussian noise process (not necessarily white), and y_n is the received signal. The top branch of the model induces multiplicative interference onto the transmitted signal point and the ISI channel introduces only the ISI. Transmission delay in the channel may be neglected

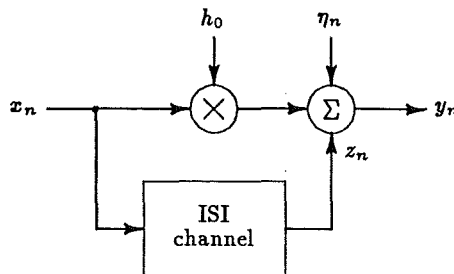


Figure 6.1 Discrete-time model of the time-dispersive channel.

and all signals are represented by equivalent lowpass signals. This channel model accounts for the combined effects of transmit and receive filters, channel characteristics, phase and timing recovery errors, and equalization.

The impulse response of the ISI channel is generally infinite in extent, but can be truncated, in practice, and modelled by a finite tapped delay line of length $k_1 + k_2 + 1$ as shown in Figure 6.2, where k_1 is the number of precursor samples and k_2 is the number of postcursor samples in the impulse response. Due to the memory elements (time delays) in the ISI channel model, channels that introduce ISI are often referred to as channels with *memory*. The time advances (D^{-1}) are included for notational convenience. The length of the delay line is chosen such that a high percentage of the peak distortion in the actual ISI channel response is accounted for in the truncated response, where peak distortion is now defined as

$$D_p \triangleq \sum_{i=-\infty}^{\infty} \left\{ |\operatorname{Re}[h_i]| + |\operatorname{Im}[h_i]| \right\} \quad (6.2)$$

The ' on the summation indicates that the $i = 0$ term is not included. The ISI in the n^{th} received symbol is given by

$$z_n = \sum_{i=-k_1}^{k_2} x_{n-i} h_i \quad (6.3)$$

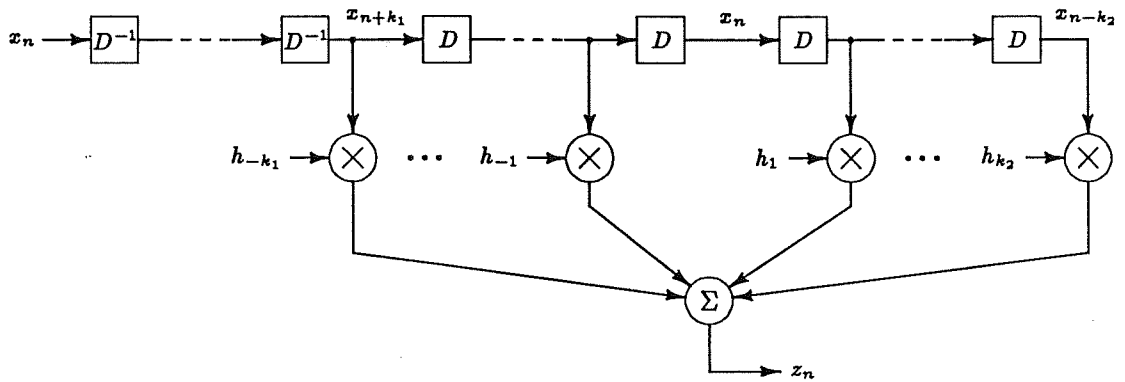


Figure 6.2 Discrete-time model of the ISI channel.

The received symbol at time n is

$$y_n = h_0 x_n + z_n + \eta_n \quad (6.4)$$

where h_0 is the multiplicative interference. Since we will only be concerned with stationary channel responses and stationary transmitted signal statistics, the pdf of the random variable Z_n , $p(z_n)$, is independent of the time n . We may therefore write it as

$$p(z) = p \left(\sum_{i=-k_1}^{k_2} x_{-i} h_i \right) \quad (6.5)$$

An exact computation of $p(z)$ generally requires knowledge of the joint pdf's of the random variables $\{X_{-i} h_i\}$ and involves $O((k_1 + k_2)2^{m(k_1 + k_2)})$ operations, where m is the number of data bits per symbol. In general, an exact computation is only possible when m and $(k_1 + k_2)$ are both small. A technique is now presented for computing an approximate ISI pdf for uncoded systems when m and $(k_1 + k_2)$ are not small.

6.2 Uncoded Systems

Consider an uncoded system where all transmitted signal points $\{x_n\}$ are mutually independent. This allows (6.5) to be expressed as the convolution

$$p(z) = p(x_{k_1} h_{-k_1}) * \cdots * p(x_1 h_{-1}) * p(x_{-1} h_1) * \cdots * p(x_{-k_2} h_{k_2}) \quad (6.6)$$

where each of the *partial pdf's* $p(x_{-i} h_i)$, is a distinct pdf and should strictly be denoted as $p_{X_{-i} h_i}(z)$, although we will continue with the above notation for convenience. This equation can be computed with $O(2^{m(k_1+k_2)})$ operations—a factor of $(k_1 + k_2)$ improvement over (6.5). However, an exact computation is still not feasible for most cases of practical interest.

The key to reducing the computational cost significantly is to quantize the ISI so that the number of distinct ISI levels is greatly reduced, and to compute an approximation to $p(z)$ using the procedure illustrated in Figure 6.3 for an uncoded system. The transmitted signal points x_{-i} are scaled by complex ISI channel impulse response samples h_i , corresponding to a scaling and a rotation of the mass points in the complex plane. Each of the partial pdf's is then binned according to a quantization level (bin width) Δ chosen to significantly reduce computational cost while incurring a small error in the approximation to $p(z)$. All mass points within a bin are summed to obtain a total mass for the bin. The total mass for each bin is then placed at the centre of the bin to obtain a binned partial pdf $\hat{p}(x_{-i} h_i)$. This process of central binning is a particularly good approximation for QAM (lattice based) constellations because the signal points are uniformly distributed across the signal space. Finally, the binned partial pdf's are convolved to obtain the approximation $\hat{p}(z)$ to $p(z)$.

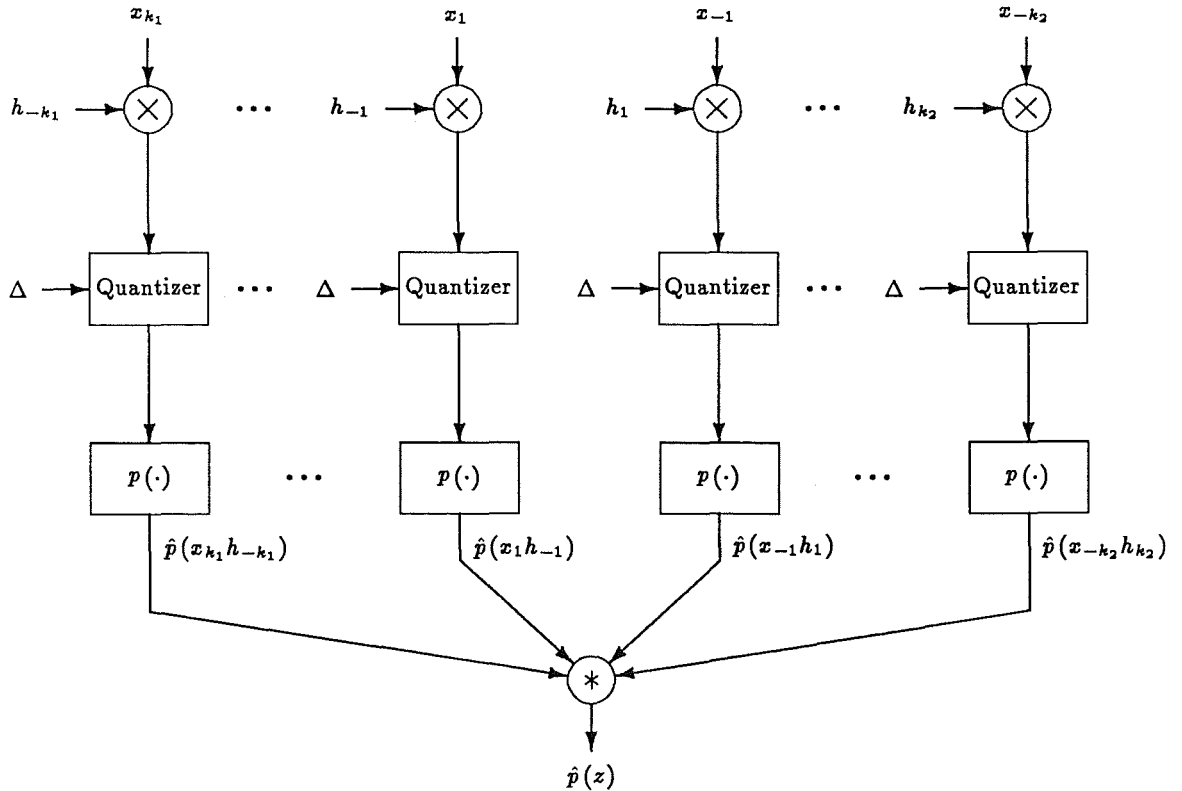
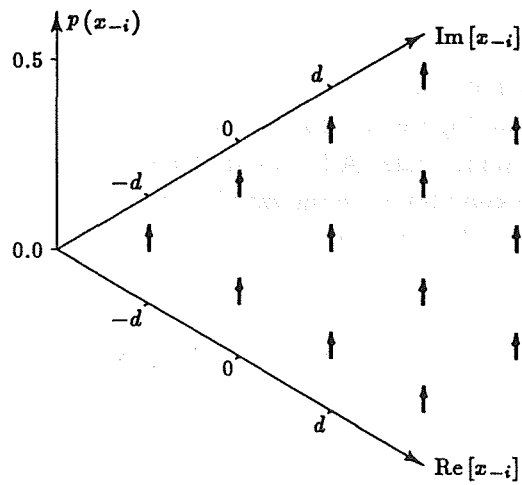
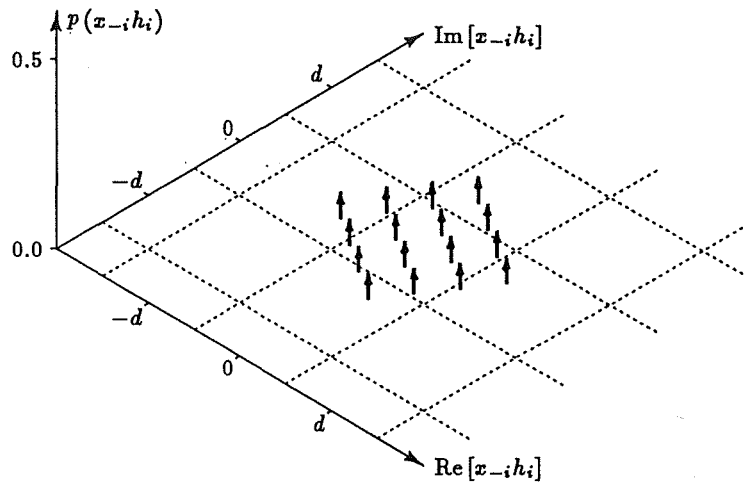


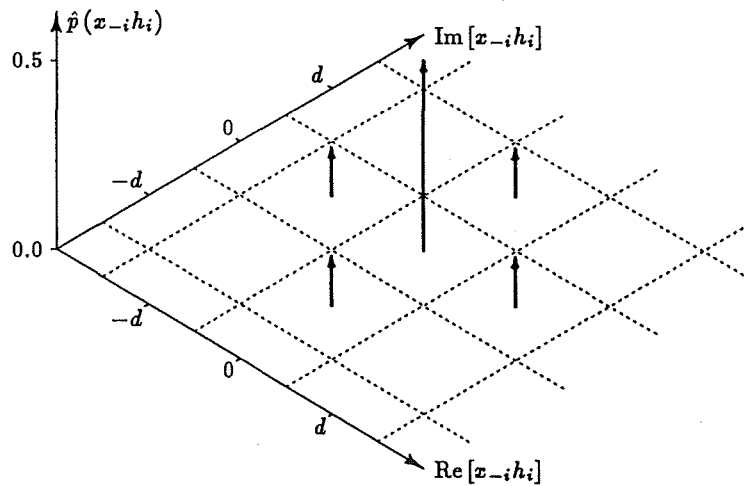
Figure 6.3 Computation of an ISI pdf for an uncoded system.



(a) Transmitted signal point pdf



(b) Partial pdf



(c) Binned partial pdf

Figure 6.4 Example pdf's during the computation of an ISI pdf ($h_i = 0.2 + j0.3$).

The procedure used to control the computational cost involves first determining the number of bins that can be easily handled in a convolution without excessive cost. A constant bin width is then chosen for all the partial pdf's so that the selected number of bins will cover the extent of the ISI pdf $p(z)$.

Some typical pdf's are shown in Figure 6.4 for the various stages of computation outlined above; they are all on the same scale. All possible transmitted signal points are assumed to be equally likely, so the complex random variable X_{-i} has a uniform discrete pdf $p(x_{-i})$ as shown in Figure 6.4a. This pdf is scaled by $h_i = 0.2 + j0.3$ to obtain a partial pdf $p(x_{-i}h_i)$ as shown in Figure 6.4b. The binned partial pdf $\hat{p}(x_{-i}h_i)$ is shown in Figure 6.4c, where the bins are denoted by the dotted lines.

There are a number of possible variations on this basic algorithm:

1. Rather than use a fixed quantization level throughout the computation, the truncated impulse response of the ISI channel can be rank ordered and the number of bins held constant for all partial pdf's. The consequence of this is that the partial pdf's of smallest extent are sampled more finely, leading to increased accuracy in the approximate ISI pdf. This technique requires $O(k_1 + k_2)$ times more computation than the basic algorithm.
2. Instead of neglecting the precursor and postcursor samples that fall outside the extent of the ISI channel impulse response, they can be lumped at respective ends of the truncated impulse response and treated as two additional samples. Alternatively, the ISI due to the truncated samples can be treated as a Gaussian random variable [Hill, 1971]. However, these techniques are not usually required if the truncated impulse response includes a sufficiently high percentage of the distortion in the infinite impulse response.
3. If the signal constellation is 90° or 180° rotationally symmetric, only $1/4$ or $1/2$ of the ISI pdf need be stored during computation because the pdf will have the same rotational symmetry.
4. The constellation decomposition previously mentioned, and used by Hill [1971], can be applied if the signal constellation is 2^m -ary PAM or 2^{2m} -ary QAM.

Normally (6.6) could be evaluated most efficiently using *fast Fourier transforms* [Oppenheim and Schaffer, 1975]. The binning of the partial pdf's, however, is a nonlinear operation that has no simple analog in the Fourier transform domain. Therefore, binning would have to be performed prior to the Fourier transforms, but this turns out to be inefficient when more than two pdf's must be convolved. Hence we resort to using discrete convolutions.

When a system incorporates channel coding, the transmitted signal points are mutually dependent, and the procedure to approximate $p(z)$ must be appropriately modified. This modification is discussed for trellis-coded systems in the next section.

6.3 Trellis-Coded Systems

The ISI pdf for a trellis-coded system cannot be expressed as in (6.6) because the transmitted signal points $\{x_n\}$ are now mutually dependent. An approximate ISI pdf cannot, therefore, be computed using the procedure described for an uncoded system. However, some of the features of TCM can be exploited to yield an efficient algorithm.

Figure 6.5 shows a decomposition that can be performed on trellis-coded constellations. Codewords have been assigned to the signal points of a 16-QAM signal set using Ungerboeck's mapping by set partitioning, with two coded and two uncoded bits. Notice that the 16-QAM signal set can be partitioned into four replicas of a *coded signal set* (CSS), where each of the replicas happen to lie in one of the four quadrants of the 16-QAM. Points in the CSS are labelled with the coded bits, and the uncoded bits can be assigned so that they are all equal for a given replica of the CSS. This means that each signal point x_n can be expressed as the *Euclidean* sum of an uncoded signal point x_n^u and a coded signal point x_n^c , as shown in Figure 6.5. The pdf of the coded signal point can therefore be expressed as a convolution of the CSS pdf, which depends on the coded bits, and the *uncoded signal set* (USS) pdf, which depends only on uncoded bits. For constellations with non-square boundaries, it may be necessary to include one or more uncoded bits in the CSS to preserve the boundary geometry in the decomposition. A similar decomposition can be performed for phase shift keying (PSK) constellations in phase space. The concept of a coded signal set is similar to the concept of a *basic signal set* as discussed by Pottie and Taylor [1987].

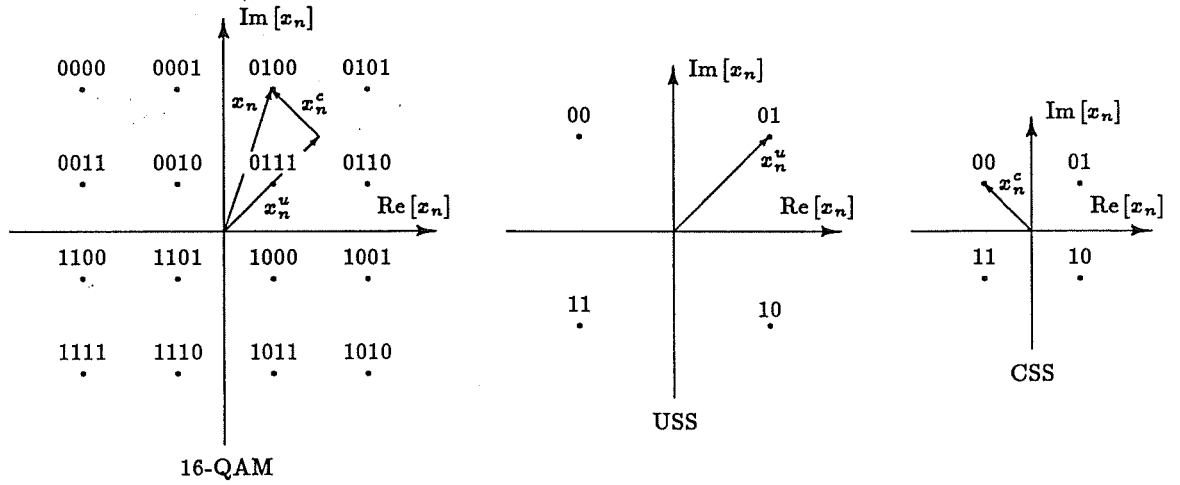


Figure 6.5 Decomposition of coded 16-QAM.

The ISI for a trellis-coded constellation can now be expressed as

$$z = \sum_{i=-k_1}^{k_2} (x_{-i}^u h_i + x_{-i}^c h_i) \quad (6.7)$$

and, because the uncoded signals are independent, the ISI pdf is given by

$$p(z) = p(x_{k_1}^u h_{-k_1}) * \dots * p(x_1^u h_{-1}) * p(x_{-1}^u h_1) * \dots * p(x_{-k_2}^u h_{k_2}) * p\left(\sum_{i=-k_1}^{k_2} x_{-i}^c h_i\right) \quad (6.8)$$

The part of $p(z)$ that depends on the uncoded signals has an identical form to (6.6) and can be approximated using the techniques described in Section 6.2. The remainder of $p(z)$ can be expressed recursively to include the effect of dependence between transmitted signal points. To derive the relationship we take a trellis of time extent (depth) $k_1 + k_2 + 1$ and consider the ISI conditioned on the trellis paths. Consider, for example, using the four state, rate 1/2 convolutional encoder in Figure 5.2a to generate the bit labels of the

CSS for the 16-QAM constellation. The trellis for this code is shown in Figure 6.6, with the possible code states $s_k \in \mathcal{S}$ during each signalling interval arranged vertically and time progressing from left to right. The state transitions are labelled with signal points from the CSS.

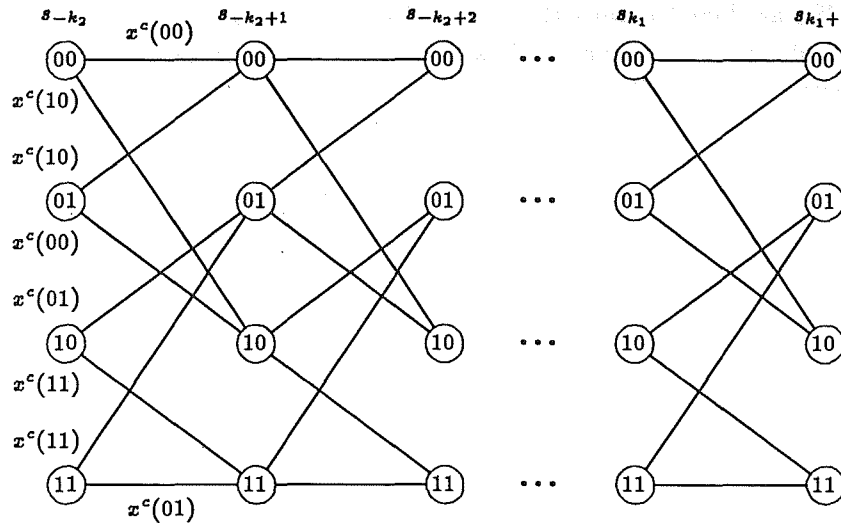


Figure 6.6 Trellis for coded 16-QAM.

We begin by conditioning on the final state transition (s_{k_1}, s_{k_1+1}) in the trellis, and work backward through the trellis. The pdf of ISI due to the coded signal can thus be expressed as

$$\begin{aligned}
 p\left(\sum_{i=-k_1}^{k_2'} x_{-i}^c h_i\right) &= \sum_{s_{k_1+1} \in \mathcal{S}} \sum_{s_{k_1} \in \mathcal{S}} p\left(\sum_{i=-k_1}^{k_2'} x_{-i}^c h_i \mid s_{k_1}, s_{k_1+1}\right) P[s_{k_1}, s_{k_1+1}] \\
 &= \sum_{s_{k_1+1} \in \mathcal{S}} \sum_{s_{k_1} \in \mathcal{S}} p\left(\sum_{i=-k_1+1}^{k_2} x_{-i}^c h_i \mid s_{k_1}\right) \\
 &\quad * p(x_{k_1}^c h_{-k_1} \mid s_{k_1}, s_{k_1+1}) \\
 &\quad \cdot P[s_{k_1}, s_{k_1+1}]
 \end{aligned} \tag{6.9}$$

where use has been made of conditional independence to introduce the convolution. The state transition probability for a code with 2^ν states and with $2^{\tilde{m}}$ states reachable from each state, assuming all state transitions to be equally likely, is given by

$$P[s_k, s_{k+1}] = \begin{cases} 1/2^{\nu+\tilde{m}} & \text{if } s'_k \in \mathcal{S}' \\ 0 & \text{if } s'_k \notin \mathcal{S}' \end{cases} \quad \text{for all } -k_2 \leq k \leq k_1 \tag{6.10}$$

where s'_k denotes the state transition (s_k, s_{k+1}) and \mathcal{S}' is the set of possible state transitions. The partial pdf in (6.9) that involves a sum of random variables can be expressed as

$$\begin{aligned}
 p\left(\sum_{i=-k_1+1}^{k_2'} x_{-i}^c h_i \mid s_{k_1}\right) &= \sum_{s_{k_1-1} \in \mathcal{S}} p\left(\sum_{i=-k_1+2}^{k_2} x_{-i}^c h_i \mid s_{k_1-1}\right) \\
 &\quad * p(x_{k_1-1}^c h_{-k_1+1} \mid s_{k_1-1}, s_{k_1}) \\
 &\quad \cdot P[s_{k_1-1} \mid s_{k_1}] \quad \text{for all } s_{k_1} \in \mathcal{S}
 \end{aligned} \tag{6.11}$$

provided $k_1 \neq 1$. If $k_1 = 1$, then

$$p \left(\sum_{i=0}^{k_2} x_{-i}^c h_i \mid s_1 \right) = \sum_{s_0 \in \mathcal{S}} p \left(\sum_{i=1}^{k_2} x_{-i}^c h_i \mid s_0 \right) \cdot P[s_0 \mid s_1] \quad \text{for all } s_1 \in \mathcal{S} \quad (6.12)$$

so that the multiplicative interference is not included in the ISI. The conditional state probability is

$$P[s_{k-1} \mid s_k] = \begin{cases} 1/2^m & \text{if } s'_k \in \mathcal{S}' \\ 0 & \text{if } s'_k \notin \mathcal{S}' \end{cases} \quad \text{for all } -k_2 \leq k \leq k_1 \quad (6.13)$$

An examination of (6.9) and (6.11) reveals that (6.9) can be computed recursively. The recursion amounts to a computation on the code trellis—similar to the Viterbi Algorithm—but using pdf's rather than metrics and not discarding any paths. After each recursion, the partial pdf's (one for each code state) are binned according to a specified ISI quantization level, in a similar fashion to the uncoded algorithm. The recursion is started with

$$p(x_{-k_2}^c h_{k_2} \mid s_{-k_2+1}) = \sum_{s_{-k_2} \in \mathcal{S}} p(x_{-k_2}^c h_{k_2} \mid s_{-k_2}, s_{-k_2+1}) \cdot P[s_{-k_2} \mid s_{-k_2+1}] \quad \text{for all } s_{-k_2+1} \in \mathcal{S} \quad (6.14)$$

For $k_2 - 1 \geq k > 0$ or $0 > k \geq -k_1$ the recursion proceeds according to

$$p \left(\sum_{i=k}^{k_2} x_{-i}^c h_i \mid s_{-k+1} \right) = \sum_{s_{-k} \in \mathcal{S}} p \left(\sum_{i=k+1}^{k_2} x_{-i}^c h_i \mid s_{-k} \right) \cdot p(x_{-k}^c h_k \mid s_{-k}, s_{-k+1}) \cdot P[s_{-k} \mid s_{-k+1}] \quad \text{for all } s_{-k+1} \in \mathcal{S} \quad (6.15)$$

If $k = 0$ then $p \left(\sum_{i=k}^{k_2} x_{-i}^c h_i \mid s_{-k+1} \right)$ is computed from (6.12). The recursion is terminated with

$$p \left(\sum_{i=-k_1}^{k_2} x_{-i}^c h_i \right) = \sum_{s_{k_1+1} \in \mathcal{S}} p \left(\sum_{i=-k_1}^{k_2} x_{-i}^c h_i \mid s_{k_1+1} \right) P[s_{k_1+1}] \quad (6.16)$$

where the state probability is

$$P[s_k] = 1/2^\nu \quad \text{for all } -k_2 \leq k \leq k_1 + 1 \quad (6.17)$$

The uncoded binned partial pdf's $\hat{p}(x_{-i}^u h_i)$, $i = -k_1, \dots, k_2$, and the coded binned partial pdf $\hat{p} \left(\sum_{i=-k_1}^{k_2} x_{-i}^c h_i \right)$ can then be convolved to compute $\hat{p}(z)$.

The technique just described for Ungerboeck codes can also be applied to other trellis codes. However, decomposition of the coded constellation into coded and uncoded sets may not always be possible, in which case an increase in computational cost must be tolerated.

6.4 Best and Worse Case Binning

The binning procedures discussed for computing approximate ISI pdf's involve mapping all samples in a bin to the centre of the bin. This is known as *central* binning and is

illustrated in Figure 6.7. Binning procedures that can be used to compute ISI pdf's that lead to upper and lower bounds on symbol error probability have also been investigated. Although this problem is difficult to study analytically, intuitively the *best case* and *worst case* binning procedures shown in Figure 6.7 should lead to best and worst case pdf's that will respectively provide lower and upper bounds on symbol error probability. Best case binning involves mapping all samples in a bin to a position that contributes the least severe ISI, whereas worst case binning involves mapping all samples in a bin to a position that contributes the most severe ISI.

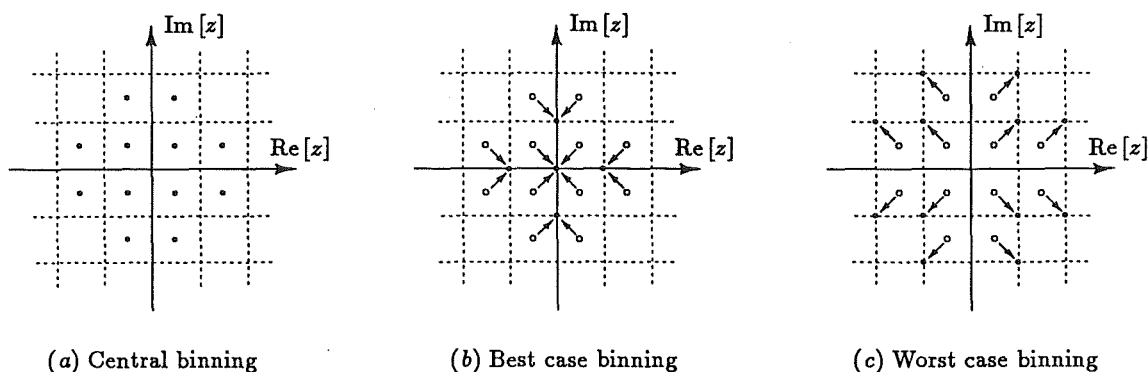


Figure 6.7 Binning procedures.

A proof that worst case binning leads to an upper bound on symbol error probability for uncoded modulation is given in Appendix 6A for a restricted but common class of ISI pdf's. From the work of Glave [1972] and Matthews [1973] it appears that tighter bounds could be obtained, but these would involve placing the mass points at arbitrary positions in the bins and would lead to a non-uniformly sampled pdf. The computational cost would then increase exponentially with $(k_1 + k_2)$, rather than linearly.

6.5 Examples of Probability Density Functions

The accuracy and utility of the described algorithms are now illustrated with some examples. These examples assume systems that incorporate raised-cosine Nyquist filtering with a roll-off factor of 0.3 and a Nyquist bandwidth of 17.5 MHz, and channels subjected to two-path fading. The impulse response of the raw ISI channel is generated by convolving the impulse response of the Nyquist filter with that of the fade, and removing the cursor from the result. The impulse response of the residual ISI channel is generated in a similar fashion, but with the addition of a five-tap synchronously-spaced equalizer, designed for minimum mean square error at high SNRs. The transfer function of the fading channel is modelled using the Rummmler model (see Section 4.2.1), with a delay $\tau = 6.3$ ns between the two paths. The signal points in a given constellation are assumed to be equally likely, so that $p(x_{-i})$ is a uniform pdf.

In the QAM constellations, the minimum distance between signal points is $d = 2$, and in the PSK constellations, all the signal points have a magnitude of one. These values have been used in the computation of the ISI levels and bin widths. A peak ISI level of one or greater causes closure of the eye.

The truncated lengths of the impulse responses of the ISI channel were chosen to account for about 99% of the distortion in the actual response; although, in general, accounting for over 90% of the distortion yielded close approximations to the actual ISI probability densities. Uncoded partial pdf's were computed using variation 1 of the basic

algorithm described in Section 6.2. For clarity of presentation, the discrete ISI pdf's were converted to continuous pdf's by linear interpolation and scaling by the bin width. Because symmetries in the signal constellations meant that all the pdf's computed were circularly symmetric, only marginal pdf's for positive $\text{Re}[z]$ are shown.

Two sets of examples are presented. The first set of examples examines the effects of bin width and binning type on the ISI pdf's. An example is also given showing that the best and worst case binnings can be used to compute bounds on symbol error probability. The second set of examples examines the form of the ISI pdf for various raw and residual ISI channels. The effect of constellation size is also examined.

Extensive comparisons of Ungerboeck coded and uncoded constellations of the same size have shown that the dependence between symbols, introduced by Ungerboeck codes, has essentially no effect on the ISI pdf. It is, however, expected that trellis codes that introduce spectral nulls [Calderbank *et al.*, 1988] will have a discernible effect on the pdf of ISI, but such codes do not feature in this thesis. In all the examples given here, only Ungerboeck coded and uncoded constellations have been considered.

6.5.1 Binning Parameters

The accuracy of an ISI pdf obtained using central binning can be studied by comparing it to the ISI pdf's obtained using best and worst case binning. The ISI pdf's from the three types of binning should converge as the bin width is reduced. A further application for ISI pdf's obtained using best and worst case binning is that they can be used to compute bounds on symbol error probability.

Examples of bound pdf's for the coded 16-QAM system discussed in Section 6.3 are given in Figure 6.8. The bound pdf's obtained by best and worst case binning and the corresponding pdf's obtained by central binning are shown in Figure 6.8a for raw ISI and

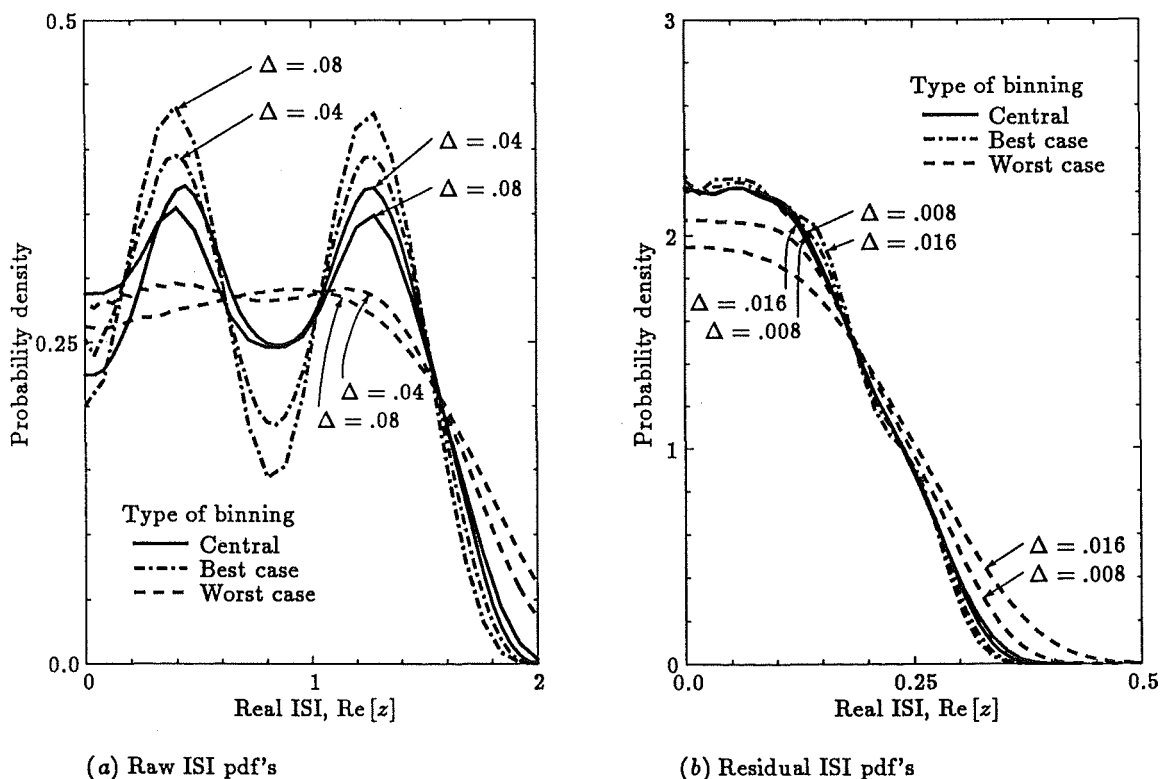


Figure 6.8 Bound pdf's of ISI for uncoded 16-QAM with $B = 20$ dB, $f_o = 0$ MHz.

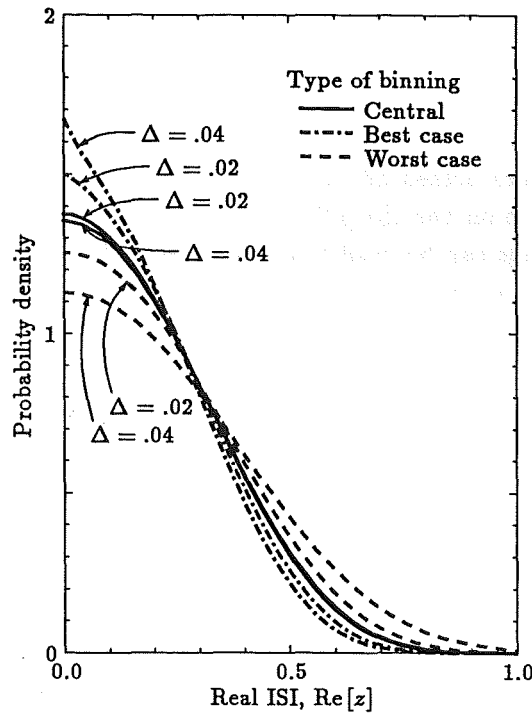


Figure 6.9 Bound pdf's of residual ISI for coded 512-CR with $B = 10$ dB, $f_o = 0$ MHz.

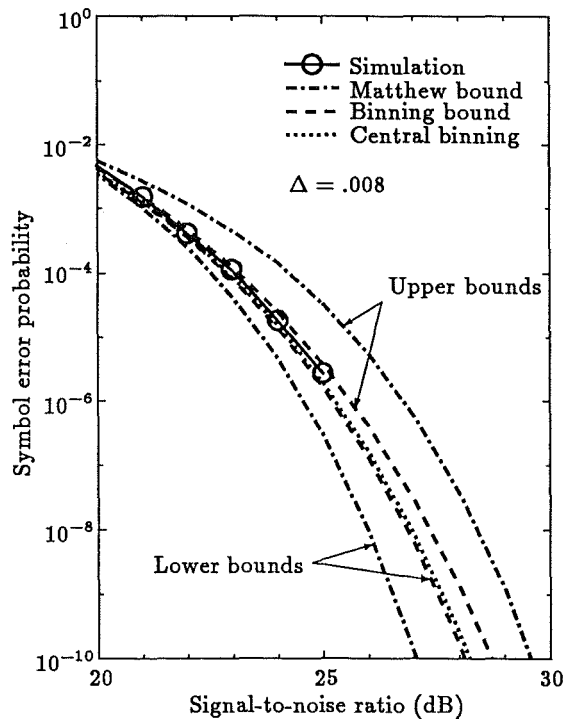


Figure 6.10 Error Bounds for uncoded 16-QAM with residual ISI with $B = 20$ dB, $f_o = 0$ MHz.

in Figure 6.8b for residual ISI. As the bin width Δ is reduced, the best and worst case pdf's approach the centrally binned pdf, which is only slightly affected by the reduced bin width. These observations suggest that the centrally binned pdf's are good approximations to the true pdf's. This was verified by performing Monte Carlo simulations to generate the ISI pdf's. The pdf's simulated from 10^6 trials are virtually indistinguishable from the centrally binned pdf's.

Bound pdf's were also computed for the $\nu = 3$ coded 512-CR constellation described in Section 5.1. The bound pdf's are shown in Figure 6.9 for two bin widths. The bell-shape of the pdf's, due to the combined effects of residual ISI and a large constellation, will be elaborated on in the following section.

Symbol error probabilities were computed as a function of signal-to-noise ratio for the uncoded 16-QAM by averaging the symbol error probability, conditioned on the ISI, over the residual ISI pdf bounds in Figure 6.8*b* (the conditional symbol error probability is similar to that for M -PAM in (6.18), but was modified for the two-dimensional QAM constellation). These symbol error probabilities are shown in Figure 6.10; the best and worst case binned pdf's respectively yield lower and upper bounds on the symbol error probability, and the centrally binned pdf yields an estimate. A curve obtained by Monte Carlo simulation is also shown, so that the tightness of the bounds can be verified. The bounds described by Matthews [1973] were modified for application to QAM constellations and computed for a comparison. The error bounds from best and worst case binnings are much tighter than Matthews' bounds for the bin width ($\Delta = .008$) used, and can be made arbitrarily tight by reducing the bin width. However, Matthews' bounds are based on bound pdf's that are optimal for the two to four mass points used, whereas the best and worst case pdf's used about one hundred mass points, took longer to compute, and are not optimal for the number of mass points.

6.5.2 Characteristics of the Probability Density Functions

Many researchers avoid computing ISI pdf's by assuming they conform to pdf's for which there are simple closed-form expressions. Typically, they model the ISI as either a Gaussian or a uniform random variable. We use the computational techniques described in this chapter to show that these are not always good approximations. For these examples, based largely on the coded 16-QAM system in Section 6.3, the bin width was chosen to be between .02 and .03 of the peak ISI. This gives a good compromise between accuracy of the ISI pdf and computational cost. All examples were computed using central binning.

In Figure 6.11*a*, the pdf's of the raw ISI are shown for various fade depths with constant delay and notch frequency; while, in Figure 6.11*b*, the notch frequency is varied for constant delay and fade depth. The latter figure shows that the raw ISI pdf's are quite

B	f_o	h_o	
		unequalized	equalized
5	4	$0.99 - j0.11$	1
10		$0.95 - j0.25$	1
15		$0.86 - j0.35$	1
20		$0.76 - j0.38$	0.99
10	0	0.99	1
	4	$0.95 - j0.25$	1
	8	$0.99 - j0.44$	1
	12	$0.99 - j0.56$	1
	16	$0.99 - j0.63$	1

Table 6.1 Multiplicative interference for 16-QAM with equalized and unequalized fades.

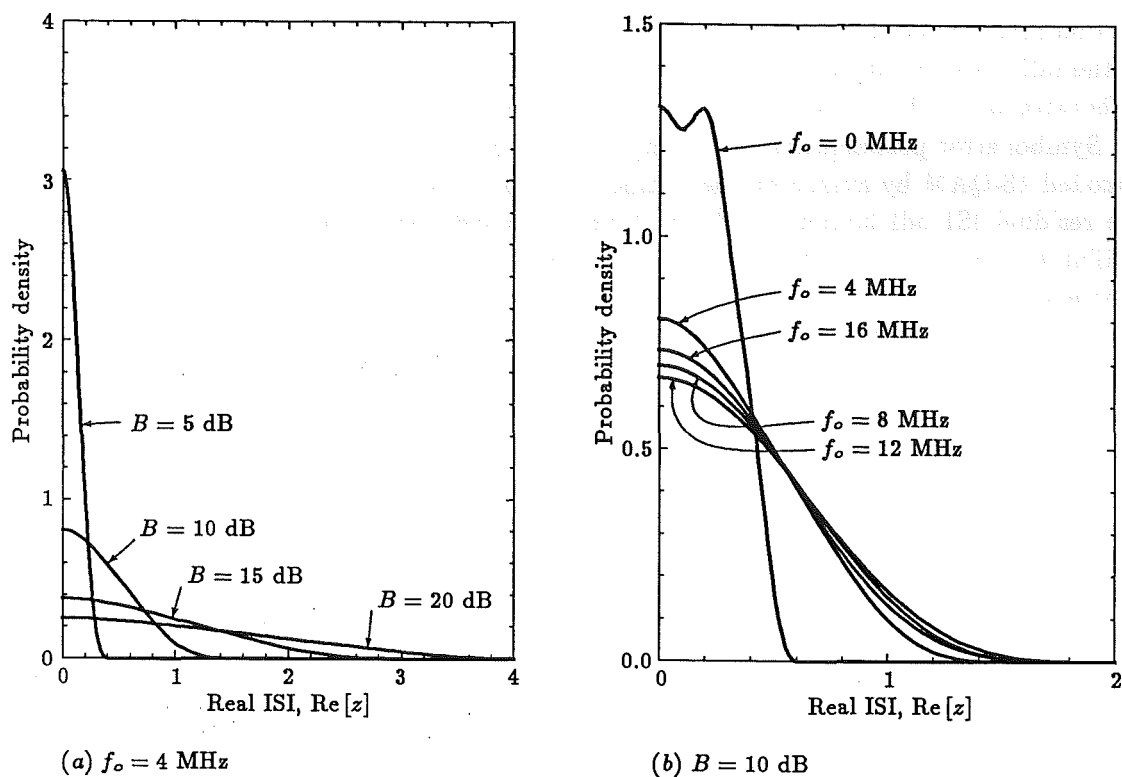


Figure 6.11 Raw ISI pdf's for coded 16-QAM.

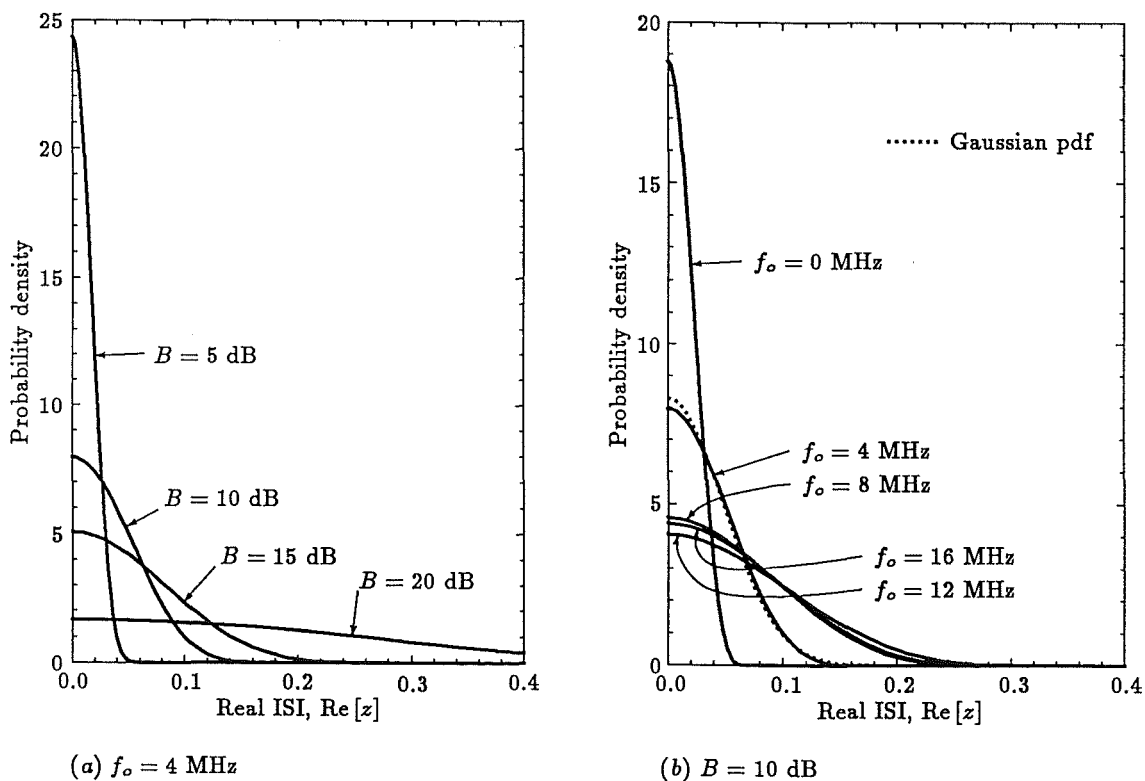
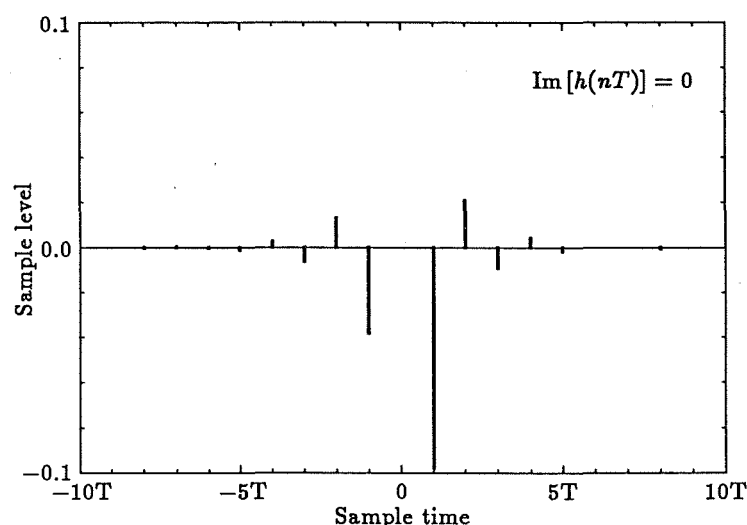


Figure 6.12 Residual ISI pdf's for coded 16-QAM.

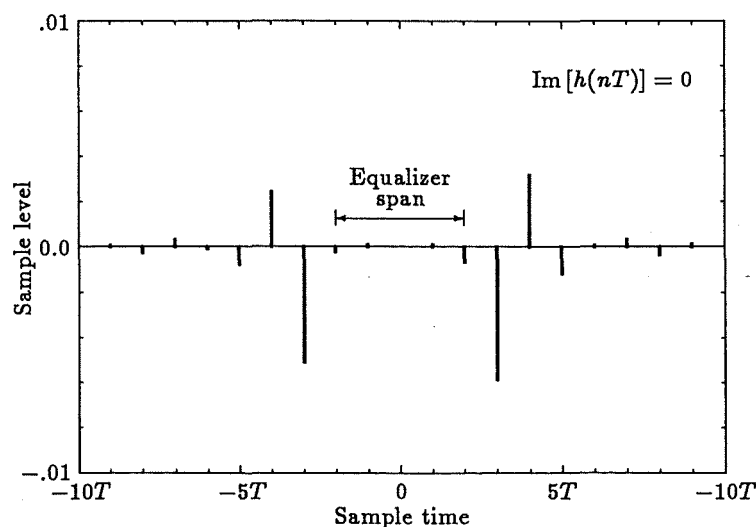
similar for the notch frequencies above 4 MHz, and this will be reflected in an almost constant symbol error probability. Figures 6.12a and b show the corresponding pdf's of the residual ISI. The equalizer significantly reduces the severity of all the raw ISI channel impulse responses, and in this case the residual ISI pdf's are quite similar for the notch frequencies above 8 MHz. Table 6.1 contains the corresponding multiplicative interference values h_0 , which indicate that the severe cross-rail interference for the unequalized fades has been cancelled in the equalized fades.

A comparison of the ISI pdf for $f_o = 4$ MHz in Figure 6.12b to the Gaussian pdf of the same variance (dotted line) shows that, although the pdf does have a bell shape, it is not actually Gaussian. A similar conclusion is reached for $f_o = 0$ MHz. In contrast, the ISI pdf for the unequalized channel with a central notch is approximately uniform, as suggested by Moridi and Sari [1985], but at other notch frequencies the uniform pdf is not a good approximation.

Inspection of the impulse responses of the raw and residual ISI channels shown in



(a) Raw ISI channel impulse response



(b) Residual ISI channel impulse response

Figure 6.13 Impulse responses for channel with $B = 10$ dB, $f_o = 0$ MHz.

Figures 6.13a and b for $f_o = 0$ MHz reveals the cause of the above observations. Notice that the impulse response of the raw ISI channel has a dominant sample, which causes one partial pdf to predominate and the raw ISI pdf to be more uniform than bell-shaped. On the other hand, the distortion of the residual ISI channel is more evenly distributed, so no partial pdf dominates and the ISI pdf is quite bell-shaped. These observations are consistent with the central limit theorem.

Another factor that affects the shape of the ISI pdf's is constellation size. Consistent with the central limit theorem, the pdf's of raw ISI in Figure 6.14 become increasingly bell-shaped as the constellation size increases. PSK constellations are used for this comparison because they all have the same amplitude and all give about the same peak ISI. If QAM constellations are used, the points in the centre of the constellation, which are lacking with PSK constellations, generate more small levels of ISI and remove the dip observed in the centre of the ISI pdf's for PSK constellations.

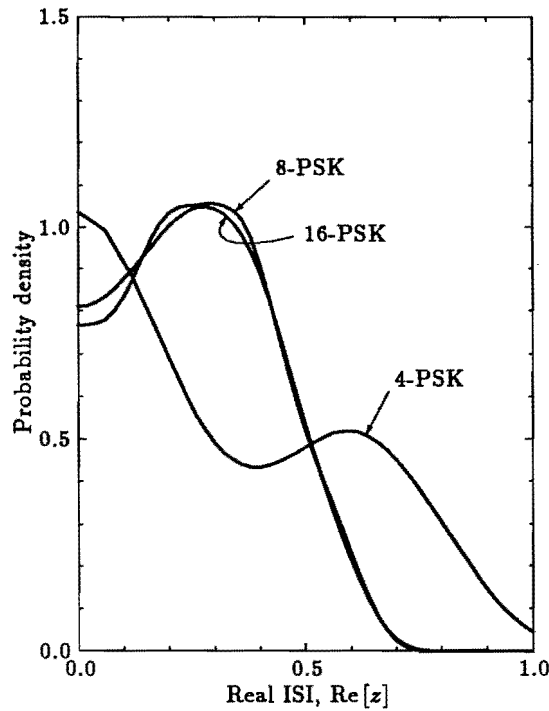


Figure 6.14 Raw ISI pdf's for various constellations with $B = 15$ dB, $f_o = 4$ MHz.

It is important to note that, although the residual ISI pdf's are quite bell shaped, ISI is peak limited and cannot in general be modelled as a Gaussian random variable. In Chapter 7, we will see that modelling the residual ISI as Gaussian can lead to extremely loose upper bounds if the ISI is severe compared to the noise.

6.6 Conclusion

An algorithm has been developed for the computation of approximate ISI probability density functions for trellis coded systems. Worst and best case binning procedures have been presented that can be used to obtain upper and lower bounds on the symbol error probability. The tightness of these bounds is controlled by the bin widths used in the computation. The error probability bounds from the worst and best case binnings are tight (provided the bin width is sufficiently small) compared to those described by Matthews,

but require a larger number of mass points and are not optimal for the number of mass points used.

The computational techniques introduced have been used to investigate the characteristics of ISI for high level QAM systems with and without trellis coding. Despite the assumptions commonly made, the pdf's of raw ISI are, in general, neither bell-shaped (due to the presence of a dominant sample in the impulse response of the ISI channel) nor uniform. Although the pdf of residual ISI is usually quite bell-shaped, it generally deviates significantly from a Gaussian pdf. The results for trellis codes of the Ungerboeck type show that the dependence between signal points in the transmitted sequence has little or no effect on the pdf of ISI.

In the next chapter, we will make use of approximate ISI pdf's, computed using the techniques in this chapter, to form bounds on decision error probability for trellis-coded modulation on time-dispersive channels.

Appendix 6A An Upper Bound on Error Probability

We prove that the worst case binning procedure leads to an upper bound on the symbol error probability for an uncoded PAM system. A similar proof exists that the best case binning procedure leads to a lower bound on error probability. Extension of these proofs to uncoded QAM systems is straightforward. The extension to trellis coded systems with Viterbi decoding is difficult due to the lack of a simple analytical expression for symbol error probability. However, it seems reasonable that the worst and best binnings would still lead to bounds on error probability.

The conditional symbol error probability for an uncoded M -PAM system with minimum distance d between signal points and ISI level z , is given by

$$P_s[\mathcal{E} | z] = \frac{M-1}{M} \left[Q\left(\frac{d-2\operatorname{Re}[z]}{2\sigma_n}\right) + Q\left(\frac{d+2\operatorname{Re}[z]}{2\sigma_n}\right) \right] \quad (6.18)$$

where $Q(\cdot)$ is the Gaussian integral function defined by (2.29), σ_n^2 is the variance of the noise, and $h_0 = 1$. The symbol error probability for this system is then

$$P_s[\mathcal{E}] = \int_{-\infty}^{\infty} P_s[\mathcal{E} | z] p(z) dz \quad (6.19)$$

$$= \int_{-\infty}^{\infty} P_s[\mathcal{E} | z] p_{-k_1}(z) * \cdots * p_{-1}(z) * p_1(z) * \cdots * p_{k_2}(z) dz \quad (6.20)$$

where $p_k(z) \equiv p_{X_{-k}h_k}(z)$ is used for notational simplicity. In general, a partial pdf can be expressed as

$$p_k(z) = \sum_{j=-J}^J P[Z = a_j] \delta(z - a_j) \quad (6.21)$$

where $P[Z = a_j]$ are the probability weights of the real-valued mass points a_j .

An unrestricted upper bound on $P_s[\mathcal{E}]$ could be obtained by applying worst case binning to the exact $p(z)$. This is intractable, however, so we apply worst case binning to each partial pdf and apply restrictions to the form of the pdf's to guarantee an upper bound on $P_s[\mathcal{E}]$. Using worst case binning procedures on the partial pdf's, the error probability is

$$P_s^w[\mathcal{E}] = \int_{-\infty}^{\infty} P_s[\mathcal{E} | z] p_{-k_1}^w(z) * \cdots * p_{-1}^w(z) * p_1^w(z) * \cdots * p_{k_2}^w(z) dz \quad (6.22)$$

where the worst case binning of a partial pdf with bin width Δ is

$$p_k^w(z) = \sum_{i=-I}^I \sum_{j \in J_i} P[Z = a_j] \delta(z - i\Delta) \quad (6.23)$$

$$= \sum_{i=-I}^I P[Z = i\Delta] \delta(z - i\Delta) \quad (6.24)$$

The set J_i defines the worst case binning

$$J_i \triangleq \begin{cases} \{j : i\Delta \leq a_j < (i+1)\Delta\} & \text{if } i < 0 \\ \{j : a_j = 0\} & \text{if } i = 0 \\ \{j : (i-1)\Delta < a_j \leq i\Delta\} & \text{if } i > 0 \end{cases} \quad (6.25)$$

We wish to prove $P_s^w[\mathcal{E}] \geq P_s[\mathcal{E}]$ for a restricted class of $p(z)$. Consider an inductive proof where we begin with the exact ISI pdf and worst case bin the partial pdf $p_k(z)$ in the k^{th} iteration for $-k_1 \leq k \leq k_2$. The pdf's *other* than the k^{th} can be expressed as

$$\Gamma_k(z) = \begin{cases} p_{-k_1+1}(z) * \cdots * p_{k_2}(z) & \text{if } k = -k_1 \\ p_{-k_1}^w(z) * \cdots * p_{k-1}^w(z) * p_{k+1}(z) * \cdots * p_{k_2}(z) & \text{if } -k_1 < k < k_2 \\ p_{-k_1}^w(z) * \cdots * p_{k_2-1}^w(z) & \text{if } k = k_2 \end{cases} \quad (6.26)$$

where $p_0(z) \equiv \delta(z)$. If we can prove

$$\int_{-\infty}^{\infty} P_s[\mathcal{E} | z] p_k(z) * \Gamma_k(z) dz \leq \int_{-\infty}^{\infty} P_s[\mathcal{E} | z] p_k^w(z) * \Gamma_k(z) dz \quad (6.27)$$

for $-k_1 \leq k \leq k_2$, then we can prove the following chain of inequalities by induction

$$\begin{aligned} & \int_{-\infty}^{\infty} P_s[\mathcal{E} | z] p_{-k_1}(z) * \cdots * p_{k_2}(z) dz \\ &= \int_{-\infty}^{\infty} P_s[\mathcal{E} | z] p_{-k_1}(z) * \Gamma_{-k_1}(z) dz \\ &\leq \int_{-\infty}^{\infty} P_s[\mathcal{E} | z] p_{-k_1}^w(z) * \Gamma_{-k_1}(z) dz \\ &= \int_{-\infty}^{\infty} P_s[\mathcal{E} | z] p_{-k_1+1}(z) * \Gamma_{-k_1+1}(z) dz \\ &\leq \int_{-\infty}^{\infty} P_s[\mathcal{E} | z] p_{-k_1+1}^w(z) * \Gamma_{-k_1+1}(z) dz \\ &\vdots \\ &\leq \int_{-\infty}^{\infty} P_s[\mathcal{E} | z] p_{k_2}^w(z) * \Gamma_{k_2}(z) dz \\ &= \int_{-\infty}^{\infty} P_s[\mathcal{E} | z] p_{-k_1}^w(z) * \cdots * p_{k_2}^w(z) dz \end{aligned} \quad (6.28)$$

and thus prove $P_s^w[\mathcal{E}] \geq P_s[\mathcal{E}]$.

In general

$$\Gamma_k(z) = \sum_{l=-L}^L P[Z = c_l] \delta(z - c_l) \quad (6.29)$$

where $P[Z = c_l]$ are the probability weights of the mass points c_l . This can be substituted into (6.27), which can be rearranged to give

$$\sum_{l=-L}^L P[c_l] \int_{-\infty}^{\infty} P_s[\mathcal{E} | z] p_k^w(z - c_l) dz \geq \sum_{l=-L}^L P[c_l] \int_{-\infty}^{\infty} P_s[\mathcal{E} | z] p_k(z - c_l) dz \quad (6.30)$$

The first restriction on the partial pdf's is that they are symmetrical about $z = 0$, which will be the case if the signal constellation is symmetrical about $x = 0$. Since $P_s[\mathcal{E} | z]$ from (6.18) is also symmetrical about zero, then to prove $P_s^w[\mathcal{E}] \geq P_s[\mathcal{E}]$ we need only prove

$$\sum_{l=-L}^L P[c_l] \int_{-\infty}^{\infty} Q\left(\frac{d - 2 \operatorname{Re}[z]}{2\sigma_n}\right) [p_k^w(z - c_l) - p_k(z - c_l)] dz \geq 0 \quad (6.31)$$

Consider the integral; using the expressions for $p_k(z)$ and $p_k^w(z)$, and the definition of J_i we get

$$\begin{aligned} & \int_{-\infty}^{\infty} Q\left(\frac{d - 2 \operatorname{Re}[z]}{2\sigma_n}\right) [p_k^w(z - c_l) - p_k(z - c_l)] dz \\ &= \int_{-\infty}^{\infty} Q\left(\frac{d - 2 \operatorname{Re}[z]}{2\sigma_n}\right) \left[\sum_{i=-I}^I P[i\Delta] \delta(z - c_l - i\Delta) \right. \\ & \quad \left. - \sum_{j=-J}^J P[a_j] \delta(z - c_l - a_j) \right] dz \end{aligned} \quad (6.32)$$

$$\begin{aligned} &= \int_{-\infty}^{\infty} Q\left(\frac{d - 2 \operatorname{Re}[z]}{2\sigma_n}\right) \left[\sum_{i=-I}^I \sum_{j \in J_i} P[a_j] \delta(z - c_l - i\Delta) \right. \\ & \quad \left. - \sum_{i=-I}^I \sum_{j \in J_i} P[a_j] \delta(z - c_l - a_j) \right] dz \end{aligned} \quad (6.33)$$

$$\begin{aligned} &= \sum_{i=1}^I \sum_{j \in J_i} P[a_j] \left[Q\left(\frac{d - 2(c_l + i\Delta)}{2\sigma_n}\right) - Q\left(\frac{d - 2(c_l + a_j)}{2\sigma_n}\right) \right. \\ & \quad \left. + Q\left(\frac{d - 2(c_l - i\Delta)}{2\sigma_n}\right) - Q\left(\frac{d - 2(c_l - a_j)}{2\sigma_n}\right) \right] \end{aligned} \quad (6.34)$$

The combination of Q functions in square brackets has the general form shown in Figure 6.15 and is positive for $c_l < d/2$, zero for $c_l = d/2$, and negative for $c_l > d/2$. Thus the integral satisfies the inequalities

$$\int_{-\infty}^{\infty} Q\left(\frac{d - 2 \operatorname{Re}[z]}{2\sigma_n}\right) [p_k^w(z - c_l) - p_k(z - c_l)] dz \begin{cases} \geq 0 & \text{if } c_l \leq d/2 \\ \leq 0 & \text{if } c_l \geq d/2 \end{cases} \quad (6.35)$$

It is now apparent that (6.31) is satisfied if $c_l \leq d/2$, which is an open eye condition on $\Gamma_k(z)$, because all the terms in the summation are positive. If the total ISI satisfies the open eye condition, then (6.31) still holds and we can prove $P_s^w[\mathcal{E}] \geq P_s[\mathcal{E}]$. When $c_l > d/2$, further restrictions must be placed on the ISI so that sum of the negative terms ($c_l > d/2$) in (6.31) does not exceed the sum of the positive terms ($c_l \leq d/2$) in magnitude. It is clear that there is a large class of closed eye conditions for which the bound also applies; however, this is difficult to define analytically because it depends on the actual form of the partial pdf's rather than just the peak ISI. Physically one restriction might be that

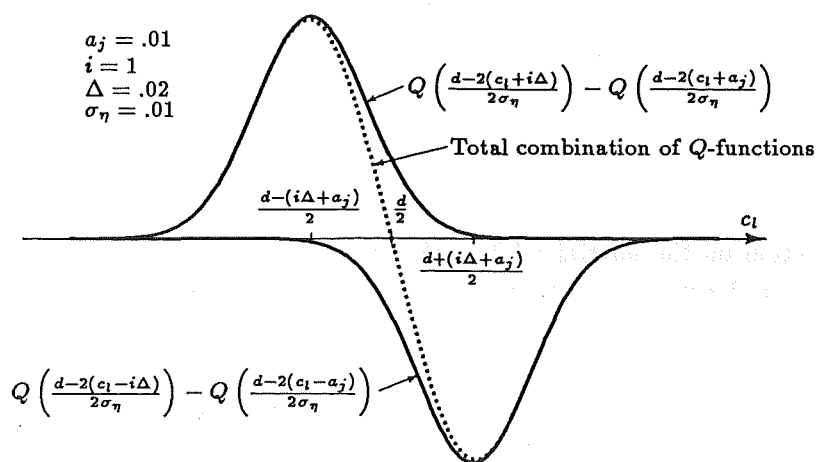


Figure 6.15 General form of the combination of Q -functions.

the tails of the ISI pdf fall away monotonically. Note that for low error probabilities ($P_s[\mathcal{E}] \leq 10^{-6}$) the eye is open with a very high probability, and it is these probabilities that we usually wish to bound analytically.

Chapter 7

Analytical Performance Bounds

There is a need for tight analytical bounds on the error probability of trellis-coded modulation (TCM) on time-dispersive channels, particularly when the Viterbi decoder only operates on the code states so that the receiver is *mismatched* to the channel. Analytical bounds can be used to estimate the performance of TCM and as tools in the design of good codes. The bounds must account for the effects of multiplicative and intersymbol interference due to receiver mismatch, which results from the time-dispersive channel, in addition to noise. We are primarily interested in bounds that account for the *residual intersymbol interference (ISI)* that exists after non-ideal equalization by a finite tap equalizer, because TCM does not perform well with *raw ISI*.

The Viterbi decoder and ISI make it difficult to obtain tractable analytical bounds on the error probability. There are, however, a number of techniques that we can apply to help solve the problem. An efficient technique to calculate upper bounds on the error event probability of a Viterbi decoder for a convolutional code was first discussed by Viterbi [1971]. This technique involves evaluating a union bound on the error-event probability using the transfer function of the convolutional code. The technique is applicable to general memoryless channels, and requires the manipulation of matrices of order $(2^\nu - 1) \times (2^\nu - 1)$ for 2^ν code states.

Bounds on the error probability of the Viterbi algorithm have also been studied by Forney [1972]. These bounds are used to analyze an optimum (matched) nonlinear receiver for uncoded data transmitted on a time-dispersive channel. The receiver consists of a whitened matched filter followed by maximum-likelihood sequence estimation using the Viterbi algorithm with a squared Euclidean distance metric. The performance of this receiver is evaluated using a union bound, an asymptotic error probability, and a lower bound.

Mismatched receivers for uncoded data on time-dispersive channels have been analyzed by Divsalar [1978]. All channel states are considered in this analysis; therefore, it is limited in practice to channels with small numbers of states. Pairwise receiver states are analyzed to evaluate a union bound. This involves the manipulation of matrices of order $(2^{2\nu} - 2^\nu) \times (2^{2\nu} - 2^\nu)$, where 2^ν is the total number of states in the receiver.

Biglieri [1984] has also considered pairwise states to evaluate a union bound for TCM on additive white Gaussian noise (AWGN) channels. The technique is applicable to general trellis codes, but requires the manipulation of matrices of order $(2^{2\nu} - 2^\nu) \times (2^{2\nu} - 2^\nu)$, where 2^ν is the number of code states. Zehavi and Wolf [1987] modified Viterbi's union bound for convolutional codes and applied it to a restricted class of linear trellis codes on AWGN channels. Their technique only requires the manipulation of matrices of order

$(2^\nu - 1) \times (2^\nu - 1)$ for 2^ν code states.

Independently of the work presented here, Oka and Biglieri [1989] have derived upper bounds on error probability for TCM on time-dispersive channels. Their technique considers all code and channel (ISI) states, so is only practical when the total number of states is small. Both matched receivers, which operate on code/channel superstates, and mismatched receivers, which only operate on the code states, are considered. This bound requires the manipulation of $(2^\nu - 1) \times (2^\nu - 1)$ matrices for 2^ν code/channel superstates. A looser upper bound, which ignores the channel states and only considers the worst possible ISI for every code path, is also presented. This bound, however, is very loose.

In this chapter we wish to analyze a mismatched receiver, but we don't want to be restricted by having to consider every channel state in addition to the code states. The bounds we derive must also be tight. We use the approximations for ISI pdf's, described in Chapter 6, and modify and extend Divsalar's technique to form useful bounds [Carlisle *et al.*, 1990c]. For 2^ν code states, we must manipulate matrices of order $(2^\nu - 1) \times (2^\nu - 1)$ when the criteria described by Zehavi and Wolf [1987] are satisfied, and matrices of order $(2^{2\nu} - 2^\nu) \times (2^{2\nu} - 2^\nu)$ otherwise.

This work is motivated by the desire to avoid simulation when studying the performance of TCM applied to digital microwave radio (DMR) systems. Simulations are time-consuming and cannot realistically provide accurate estimates of $P[\mathcal{E}]$ below 10^{-6} . This is often the region of interest for data transmission [CCITT, 1988].

The additional mathematical notation required for this chapter is presented in Section 7.1. A union bound on the error-event probability of the Viterbi decoder is formed in Section 7.2. This bound assumes that the ISI pdf, or a good approximation to it, is known. To evaluate the union bound, the pairwise error probability, conditioned on the ISI, must also be upper bounded. Several upper bounds on the conditional pairwise error probability are presented in Section 7.3. The union bound can be evaluated numerically using the techniques described in Section 7.4, and the ISI-degraded minimum distance of the TCM can be computed numerically as described in Section 7.5. A simple lower bound is given by the system error probability without ISI; this is presented in Section 7.6. Finally, examples are studied in Section 7.7 to verify the usefulness of the bounds. The upper bounds are shown to be tight for a wide range of channel conditions when compared to simulation results, and can be used for either raw or residual ISI. The lower bound is tight for low levels of ISI, but loose for severe ISI.

7.1 Preliminaries

We will use the discrete-time model of the channel introduced in Chapter 6 and illustrated in Figure 6.1. The Gaussian noise samples will usually be correlated after adaptive equalization, but the correlation is assumed to have a negligible effect on the error probability. This assumption is particularly good when the variance of the ISI is large compared to the variance of the noise. Consideration of such dependencies would unnecessarily complicate the error analysis. However, noise enhancement due to the equalizer can be significant with severe ISI and will be included in the analysis.

The discrete-time model of the system we wish to analyze, incorporating the discrete-time model of the channel, is shown in Figure 7.1 where, at time n , $u_n \in \mathcal{U}$ is the source symbol and $\hat{u}_n \in \mathcal{U}$ is the decoded source symbol. The Ungerboeck encoder can be specified by the mapping functions defined in Section 3.2.1.

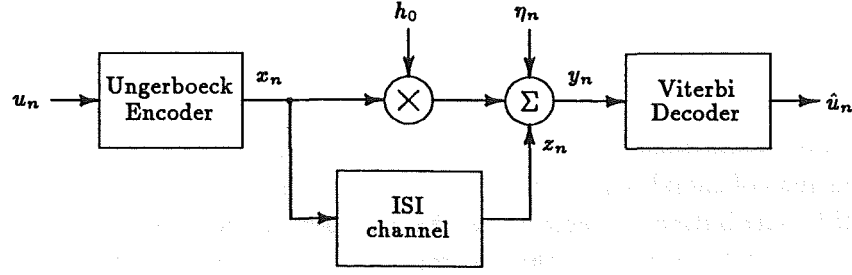


Figure 7.1 Discrete-time model of the system to be analyzed.

The following finite length sequences are used in the analysis:

$$\begin{aligned}
 \text{source symbol sequence} \quad \mathbf{u} &\triangleq (u_0, u_1, \dots, u_{N-1}) \in \mathcal{U}_N \\
 \text{state sequence} \quad \mathbf{s} &\triangleq (s_0, s_1, \dots, s_{N-1}, s_N) \in \mathcal{S}_{N+1} \\
 \text{state transition sequence} \quad \mathbf{s}' &\triangleq (s'_0, s'_1, \dots, s'_{N-1}) \in \mathcal{S}'_N \\
 \text{signal point sequence} \quad \mathbf{x} &\triangleq (x_0, x_1, \dots, x_{N-1}) \in \mathcal{X}_N \\
 \text{ISI sequence} \quad \mathbf{z} &\triangleq (z_0, z_1, \dots, z_{N-1}) \in \mathcal{Z}_N(\mathbf{x})
 \end{aligned} \tag{7.1}$$

where s'_n represents a state transition (s_n, s_{n+1}) . A hat on any of the above variables denotes the corresponding variable in the decoder, for example $\hat{\mathbf{x}} \in \mathcal{X}_N$ is the decoded signal point sequence. An error event of length $N > 1$ occurs at time $n = 0$, when the decoded state sequence diverges from and later merges with the transmitted state sequence. This is defined by restricting the possible decoded state sequences to the set

$$\mathcal{S}'_N(\mathbf{s}') \triangleq \{\hat{\mathbf{s}}' : \hat{s}_0 = s_0, \hat{s}_N = s_N, \hat{s}_n \neq s_n \text{ for } 0 < n < N\} \tag{7.2}$$

$\mathcal{Z}_N(\mathbf{x})$ is the set of ISI sequences that can disturb the transmitted signal point sequence \mathbf{x} . $\mathcal{X}_N(\mathbf{s}')$ is the set of transmitted signal point sequences for a given state transition sequence. $\mathcal{X}(\mathbf{s}'_n)$ is the set of transmitted signal points for a given state transition. The other sequence sets are unrestricted.

7.2 Union Bound on Error-Event Probability

The probability of an error event starting at time $n = 0$ is given by

$$P_e[\mathcal{E}] = P \left[\bigcup_{N=1}^{\infty} \bigcup_{\mathbf{s}' \in \mathcal{S}'_N} \bigcup_{\hat{\mathbf{s}}' \in \mathcal{S}'_N(\mathbf{s}')} \bigcup_{\mathbf{x} \in \mathcal{X}_N(\mathbf{s}')} \bigcup_{\substack{\hat{\mathbf{x}} \in \mathcal{X}_N(\hat{\mathbf{s}}') \\ \hat{\mathbf{x}} \neq \mathbf{x}}} \bigcup_{\mathbf{z} \in \mathcal{Z}_N(\mathbf{x})} \{\hat{\mathbf{x}}, \mathbf{x}, \mathbf{z}\} \right] \tag{7.3}$$

which is the probability of the union of all the error events that start at time $n = 0$. The restriction $\hat{\mathbf{x}} \neq \mathbf{x}$ is necessary to specify error events between parallel transitions ($N = 1$).

A countably infinite number of error events, which are not disjoint [Foschini, 1975], must be considered to compute $P_e[\mathcal{E}]$; therefore, (7.3) is difficult to evaluate. However, a union bound on (7.3) yields the upper bound on the error-event probability

$$P_e[\mathcal{E}] \leq \sum_{N=1}^{\infty} \sum_{\mathbf{s}' \in \mathcal{S}'_N} \sum_{\hat{\mathbf{s}}' \in \mathcal{S}'_N(\mathbf{s}')} \sum_{\mathbf{x} \in \mathcal{X}_N(\mathbf{s}')} \sum_{\substack{\hat{\mathbf{x}} \in \mathcal{X}_N(\hat{\mathbf{s}}') \\ \hat{\mathbf{x}} \neq \mathbf{x}}} \sum_{\mathbf{z} \in \mathcal{Z}_N(\mathbf{x})} P[\hat{\mathbf{x}} | \mathbf{x}, \mathbf{z}] P[\mathbf{x}] P[\mathbf{z} | \mathbf{x}] \tag{7.4}$$

where $P[\hat{\mathbf{x}} | \mathbf{x}, \mathbf{z}]$ is known as the *conditional pairwise error probability*.

The nature of TCM allows us to assume that \mathbf{z} is independent of \mathbf{x} and that the ISI samples $\{z_n\}$ are mutually independent. Two aspects of TCM can be identified that support these assumptions of independence. First, the set partitioning ensures that the parallel transitions associated with any branch of the code trellis map to subsets with similar distributions of signal points from the constellation. Thus all state sequences have very similar ISI distributions. Second, as the number of parallel transitions increases (number of uncoded bits increases), the assumptions of independence become more exact because the subsets of signal points contain more points. This has the effect of making the ISI probability distributions for all state sequences even more similar. Alternatively, we can view this effect as the uncoded (independent) portion of the system dominating over the coded (dependent) portion.

When the ISI samples $\{z_n\}$ are mutually dependent, the transmitted signal point sequence can be interleaved [Viterbi and Omura, 1979] at the transmitter and the received signal sequence can be deinterleaved at the receiver so that the samples appear independent at the receiver. If a system uses interleaving with an infinite interleaving depth, then the independence assumptions become exact. For systems without interleaving and with small constellations, the assumptions may be inaccurate. In such cases either independence can be assumed or the technique described by Oka and Biglieri [1989] can be used. The effect of interleaving will be studied in Section 7.7 using Monte Carlo simulation.

Applying the above assumptions of independence to (7.4), we get

$$P_e[\mathcal{E}] \leq \sum_{N=1}^{\infty} \sum_{s' \in \mathcal{S}'_N} \sum_{\hat{s}' \in \mathcal{S}'_N(s')} \sum_{\mathbf{x} \in \mathcal{X}_N(s')} \sum_{\substack{\hat{\mathbf{x}} \in \mathcal{X}_N(\hat{s}') \\ \hat{\mathbf{x}} \neq \mathbf{x}}} \sum_{z_0 \in \mathcal{Z}} \cdots \sum_{z_{N-1} \in \mathcal{Z}} P[\hat{\mathbf{x}} | \mathbf{x}, \mathbf{z}] P[\mathbf{x}] \cdot P[z_0] \cdots P[z_{N-1}] \quad (7.5)$$

The probability of the ISI $P[z_i]$ can be approximated using the techniques described in Chapter 6, and $P[\mathbf{x}]$ is simple to determine for Ungerboeck codes when the probability of a given source symbol $P[\mathbf{u}]$ is known. Although $P[\hat{\mathbf{x}} | \mathbf{x}, \mathbf{z}]$ can be evaluated for given \mathbf{x} , $\hat{\mathbf{x}}$, and \mathbf{z} , there is a countably infinite number of pairwise errors, so further simplification of the bound is required. If $P[\hat{\mathbf{x}} | \mathbf{x}, \mathbf{z}]$ can be expressed in the form of a product (usually of exponentials) over time n , then a generalization of the transfer function of the trellis code can be used to evaluate the bound in (7.5) without having to explicitly compute every term. Upper bounds on the conditional pairwise error probability are discussed in the next section.

7.3 Upper Bounds on Conditional Pairwise Error Probability

A Viterbi decoder selects the maximum-likelihood signal point sequence using a decoding metric. A metric that is commonly used is the squared Euclidean distance metric, which is optimum for AWGN channels, but may not be optimum for time-dispersive channels. Nevertheless, in the absence of a better metric, this is the metric we assume in this chapter. The contribution to the squared Euclidean distance metric at time n is

$$m(y_n, x_n) \triangleq -|y_n - x_n|^2 \quad (7.6)$$

A necessary and sufficient condition for an error event involving $\hat{\mathbf{x}}$ to occur at time

$n = 0$ is [Foschini, 1975]

$$\sum_{n=0}^{N-1} m(y_n, \hat{x}_n) \geq \sum_{n=0}^{N-1} m(y_n, \tilde{x}_n) \quad \text{for all } \tilde{x} \neq \hat{x}, \tilde{x} \in \mathcal{X}_N \quad (7.7)$$

provided that the condition for an error event specified by (7.2) is satisfied.

We can form an upper bound on $P[\hat{x} | \mathbf{x}, \mathbf{z}]$ by only requiring that (7.7) is satisfied for the transmitted sequence \mathbf{x} (rather than for all $\tilde{x} \neq \hat{x}$). This weakens the condition for an error event involving \hat{x} and loosens the upper bound on the pairwise error probability to give

$$P[\hat{x} | \mathbf{x}, \mathbf{z}] \leq P \left[\left\{ \sum_{n=0}^{N-1} m(y_n, \hat{x}_n) \geq \sum_{n=0}^{N-1} m(y_n, x_n) \right\} \middle| \mathbf{x}, \mathbf{z} \right] \quad (7.8)$$

If we define the metric difference

$$\Delta M \triangleq \sum_{n=0}^{N-1} \{m(y_n, \hat{x}_n) - m(y_n, x_n)\} \quad (7.9)$$

then the bound on the conditional pairwise error probability is

$$P[\hat{x} | \mathbf{x}, \mathbf{z}] \leq P[\Delta M \geq 0 | \mathbf{x}, \mathbf{z}] \quad (7.10)$$

The bound in (7.10) must now be expressed in the form of a product of terms over time n so that it can be evaluated using a generalization of the transfer function of the trellis code. There are a number of ways this can be achieved, all of which result in a weakening of the upper bound. We will examine one technique that uses a Viterbi bound and another technique that uses a Chernoff bound.

7.3.1 Viterbi Bound

The metric difference in (7.9) can be expanded by substituting for y_n from (6.4) and for $m(\cdot, \cdot)$ from (7.6) to get

$$\Delta M = \sum_{n=0}^{N-1} \left\{ -|h_0 x_n + z_n + \eta_n - \hat{x}_n|^2 + |h_0 x_n + z_n + \eta_n - x_n|^2 \right\} \quad (7.11)$$

which can be rearranged to obtain

$$\Delta M = - \sum_{n=0}^{N-1} \left\{ |x_n - \hat{x}_n|^2 + 2 \operatorname{Re} [(x_n - \hat{x}_n)^* (x_n(h_0 - 1) + z_n + \eta_n)] \right\} \quad (7.12)$$

For specific values of \mathbf{x} , $\hat{\mathbf{x}}$, and \mathbf{z} , ΔM is a Gaussian random variable with mean

$$\mu_{\Delta M} = - \sum_{n=0}^{N-1} \left\{ |x_n - \hat{x}_n|^2 + 2 \operatorname{Re} [(x_n - \hat{x}_n)^* (x_n(h_0 - 1) + z_n)] \right\} \quad (7.13)$$

and variance

$$\sigma_{\Delta M}^2 = \sum_{n=0}^{N-1} 4\sigma_{\eta}^2 |x_n - \hat{x}_n|^2 \quad (7.14)$$

where σ_{η}^2 is the variance of the Gaussian noise at the output of the equalizer (i.e. including the effect of noise enhancement). Thus the conditional pairwise error probability can be bounded by the Gaussian integral function (Q -function)

$$P[\hat{x} | \mathbf{x}, \mathbf{z}] \leq P[\Delta M \geq 0 | \mathbf{x}, \mathbf{z}] = Q \left(\frac{-\mu_{\Delta M}}{\sigma_{\Delta M}} \right) \quad (7.15)$$

where

$$Q(s) = \int_s^\infty e^{-t^2/2} dt \quad (7.16)$$

To bound the conditional pairwise error probability by a product over time n , we use the Viterbi upper bound on the Q -function (see Section 2.3.3). We apply this bound in the form

$$Q\left(\sqrt{\frac{d_{min}^2 + l^2}{4\sigma_\eta^2}}\right) \leq \exp\left(\frac{-l^2}{8\sigma_\eta^2}\right) Q\left(\frac{d_{min}}{2\sigma_\eta}\right) \quad \text{for } d_{min}^2, l^2 \geq 0 \quad (7.17)$$

where d_{min}^2 is a constant that will be computed to lower bound the ISI-degraded minimum distance of the code. It is tightest when l^2 is small compared to d_{min}^2 . In our application this corresponds to the conditional pairwise error probabilities that contribute most significantly to $P_e[\mathcal{E}]$. When $l^2 \gg d_{min}^2$, the Viterbi bound becomes an exponential bound [Wozencraft and Jacobs, 1965].

To use the Viterbi bound, we let

$$\sqrt{\frac{d_{min}^2 + l^2}{4\sigma_\eta^2}} = \frac{-\mu_{\Delta M}}{\sigma_{\Delta M}} \quad (7.18)$$

so that

$$\begin{aligned} d_{min}^2 + l^2 &= \sum_{n=0}^{N-1} \left\{ |x_n - \hat{x}_n|^2 + 4 \operatorname{Re}[(x_n - \hat{x}_n)^*(x_n(h_0 - 1) + z_n)] \right\} \\ &\quad + \frac{\left\{ \sum_{n=0}^{N-1} 2 \operatorname{Re}[(x_n - \hat{x}_n)^*(x_n(h_0 - 1) + z_n)] \right\}^2}{\sum_{n=0}^{N-1} |x_n - \hat{x}_n|^2} \end{aligned} \quad (7.19)$$

This expression must be put in the form of a summation $\sum_{n=0}^{N-1} q(x_n, \hat{x}_n, z_n)$ so that the conditional pairwise error probability can be bounded by a product over time n . Unfortunately, (7.19) cannot be expressed by a simple summation, so we require a lower bound for $d_{min}^2 + l^2$ that can be expressed by a simple summation. Divsalar [1978] uses such a bound in his work on mismatched receivers. He forms a lower bound on a/\sqrt{b} of the form

$$\frac{a}{\sqrt{b}} \geq \sqrt{4\rho(a - \rho b)} \quad (7.20)$$

for $b \geq 0$ and $a \geq \rho b$, and where $0 \leq \rho \leq 1/2$ is a parameter of the bound. For our purposes, this bound can be expressed as

$$\sqrt{d_{min}^2 + l^2} = \frac{-\mu_{\Delta M}}{\sigma_{\Delta M}/2\sigma_\eta} \geq \sqrt{-4\rho(\mu_{\Delta M} + \rho\sigma_{\Delta M}^2/4\sigma_\eta^2)} \quad (7.21)$$

This bound is valid provided $\sigma_{\Delta M}^2 \geq 0$, which it always is, and $-\mu_{\Delta M} \geq \rho\sigma_{\Delta M}^2/4\sigma_\eta^2$, which is satisfied if the ISI does not exceed some limit—to be determined. Substituting for $\mu_{\Delta M}$ and $\sigma_{\Delta M}^2$, we find that

$$d_{min}^2 + l^2 \geq 4\rho \sum_{n=0}^{N-1} \left\{ (1 - \rho)|x_n - \hat{x}_n|^2 + 2 \operatorname{Re}[(x_n - \hat{x}_n)^*(x_n(h_0 - 1) + z_n)] \right\} \quad (7.22)$$

which is a lower bound in the form of a summation. The Q -function in (7.15) is a monotonically decreasing function of $-\mu_{\Delta M}/\sigma_{\Delta M}$, so can now be upper bounded as

$$Q\left(\frac{-\mu_{\Delta M}}{\sigma_{\Delta M}}\right) \leq Q\left(\sqrt{\frac{-4\rho(\mu_{\Delta M} + \rho\sigma_{\Delta M}^2/4\sigma_\eta^2)}{4\sigma_\eta^2}}\right) \quad (7.23)$$

The Viterbi bound in (7.17) can be applied to this bound to get an upper bound that can be expressed in the form of a product of exponential terms over time n

$$Q\left(\frac{-\mu_{\Delta M}}{\sigma_{\Delta M}}\right) \leq Q\left(\frac{d_{\min}(\rho)}{2\sigma_{\eta}}\right) \exp\left(\frac{d_{\min}^2(\rho)}{8\sigma_{\eta}^2}\right) \exp\left(\frac{\rho(\mu_{\Delta M} + \rho\sigma_{\Delta M}^2/4\sigma_{\eta}^2)}{2\sigma_{\eta}^2}\right) \quad (7.24)$$

where the minimum squared distance $d_{\min}^2(\rho)$ is defined as

$$d_{\min}^2(\rho) \triangleq \min_{N, \mathbf{x}, \hat{\mathbf{x}} \neq \mathbf{x}, \mathbf{z}} \left\{ -4\rho \left(\mu_{\Delta M} + \rho\sigma_{\Delta M}^2/4\sigma_{\eta}^2 \right) \right\} \quad (7.25)$$

which is less than or equal to d_{\min}^2 in (7.17) because it is minimized over all conditional pairwise errors. Note that $d_{\min}^2(\rho) \geq 0$ provided $-\mu_{\Delta M} \geq \rho\sigma_{\Delta M}^2/4\sigma_{\eta}^2$, which is the same condition as for (7.21). If this condition is not satisfied, however, $d_{\min}^2(\rho) = 0$ because it cannot be less than zero.

The bound in (7.24) must be optimized over the ρ parameter. However, it is not feasible to optimize each conditional pairwise error probability, so we will use the same ρ for all the conditional pairwise error probability bounds and optimize the resulting union bound over ρ .

For a QAM signal constellation with a minimum distance d between the signal points, we find that the condition for which the bound in (7.21) applies (i.e. $-\mu_{\Delta M} \geq \rho\sigma_{\Delta M}^2/4\sigma_{\eta}^2$) can be used to approximately derive the limit on the ISI for which (7.21) is valid. This limit must be minimized over all conditional pairwise errors and is derived in Appendix 7A as

$$z_{\max} = (1 - \rho)d/2 \quad (7.26)$$

The value of z_{\max} is $d/2$ when $\rho = 0$; this corresponds to an eye closure threshold. When the ISI exceeds z_{\max} , the bound in (7.24) still applies, provided that $d_{\min}^2(\rho) = 0$, where d_{\min} is the minimum possible minimum distance. However, the bound exceeds $1/2$, so it becomes very loose. In practice, there will often only be a small probability that the ISI exceeds the limit, so although some of the conditional pairwise error probability bounds will be loose, the union bound can still be tight.

In the absence of ISI, the optimum bound parameter value is $\rho = 1/2$. This value can also be used when the ISI is small, to avoid having to optimize the bound. For $\rho = 1/2$, the bound on the Q -function becomes

$$Q\left(\frac{-\mu_{\Delta M}}{\sigma_{\Delta M}}\right) \leq Q\left(\frac{d_{\min}(1/2)}{2\sigma_{\eta}}\right) \exp\left(\frac{d_{\min}^2(1/2)}{8\sigma_{\eta}^2}\right) \cdot \exp\left(-\frac{\sum_{n=0}^{N-1} \{|x_n - \hat{x}_n|^2 + 4 \operatorname{Re}[(x_n - \hat{x}_n)^*(x_n(h_0 - 1) + z_n)]\}}{8\sigma_{\eta}^2}\right) \quad (7.27)$$

where $d_{\min}^2(1/2)$ is defined as

$$d_{\min}^2(1/2) \triangleq \min_{N, \mathbf{x}, \hat{\mathbf{x}} \neq \mathbf{x}, \mathbf{z}} \left\{ \sum_{n=0}^{N-1} \left\{ |x_n - \hat{x}_n|^2 + 4 \operatorname{Re}[(x_n - \hat{x}_n)^*(x_n(h_0 - 1) + z_n)] \right\} \right\} \quad (7.28)$$

The limit on the ISI for $\rho = 1/2$ is given by (7.26) as $z_{\max} = d/4$.

7.3.2 Chernoff Bound

An upper bound on the conditional pairwise error probability can also be formed by applying a Chernoff bound (see Section 2.3.3) to the bound in (7.10). The Chernoff

bound is looser than the Viterbi bound for small levels of ISI, but we will see that it can be tighter for severe ISI. We also expect the bound to be tighter than the Saltzberg bound [Saltzberg, 1968], which uses a Chernoff bound over the noise and ISI, but partitions the channel impulse response in two and treats the ISI due to one part as Gaussian and the ISI due to the other part in a worst case fashion. A Saltzberg bound could be used here to avoid computing an approximation to the ISI pdf, but the bound would be looser because the approximation of the ISI pdf is less accurate.

A Chernoff bound on the random variable ΔM in (7.10) gives

$$\begin{aligned} P[\Delta M \geq 0 \mid \mathbf{x}, \mathbf{z}] &\leq E[\exp(\lambda \Delta M) \mid \mathbf{x}, \mathbf{z}] \quad \text{for } \lambda \geq 0 \\ &\leq \prod_{n=0}^{N-1} \exp(-\lambda |x_n - \hat{x}_n|^2) \\ &\quad \cdot \exp(-2\lambda \operatorname{Re}[(x_n - \hat{x}_n)^*(x_n(h_0 - 1) + z_n)]) \\ &\quad \cdot E[\exp(-2\lambda \operatorname{Re}[(x_n - \hat{x}_n)^*\eta_n]) \mid \mathbf{x}, \mathbf{z}] \end{aligned} \quad (7.29)$$

where the expectation is over the noise samples, which are independent of \mathbf{x} and \mathbf{z} , and λ is the bound parameter. The noise samples are also assumed to be mutually independent (i.e. the Gaussian noise is assumed to be uncorrelated). Assuming $\operatorname{Re}[\eta_n]$ and $\operatorname{Im}[\eta_n]$ are uncorrelated, the expectation

$$E[\exp(-2\lambda \operatorname{Re}[(x_n - \hat{x}_n)^*\eta_n]) \mid \mathbf{x}, \mathbf{z}] = \exp(2\lambda^2 \sigma_\eta^2 |x_n - \hat{x}_n|^2) \quad (7.30)$$

and can be used to obtain

$$\begin{aligned} P[\Delta M \geq 0 \mid \mathbf{x}, \mathbf{z}] &\leq \prod_{n=0}^{N-1} \exp(-\lambda |x_n - \hat{x}_n|^2 (1 - 2\lambda \sigma_\eta^2)) \\ &\quad \cdot \exp(-2\lambda \operatorname{Re}[(x_n - \hat{x}_n)^*(x_n(h_0 - 1) + z_n)]) \end{aligned} \quad (7.31)$$

The Chernoff parameter λ to optimize the tightness of (7.31) must satisfy

$$\frac{\partial}{\partial \lambda} E[\exp(\lambda \Delta M) \mid \mathbf{x}, \mathbf{z}] = 0 \quad (7.32)$$

However, λ depends on \mathbf{z} so it must be averaged over the ISI. If the ISI is small compared to the noise ($\sigma_z^2 < \sigma_\eta^2$), then a value of λ that is nearly optimum can be obtained by assuming the ISI has a Gaussian pdf, to get

$$\lambda = \frac{1}{4(\sigma_\eta^2 + \sigma_z^2)} \quad (7.33)$$

However, when $\sigma_z^2 \geq \sigma_\eta^2$ this is a bad assumption because the ISI is peak limited, whereas a Gaussian random variable is not. In this case (7.31) must first be averaged over the ISI and then optimized numerically to minimize the pairwise error probability bound for each value of σ_η^2 . As with the Viterbi bound, it is not feasible to optimize each conditional pairwise error probability bound, so we use the same λ for all the conditional pairwise error probability bounds and use it to optimize the union bound.

7.4 Numerical Evaluation of the Union Bound

Using either of the techniques described in the previous section, the conditional pairwise error probability can be upper bounded by a product

$$P[\hat{\mathbf{x}} \mid \mathbf{x}, \mathbf{z}] \leq \prod_{n=0}^{N-1} e^{q(x_n, \hat{x}_n, z_n)} \quad (7.34)$$

where, for the Viterbi bound, $q(\cdot, \cdot, \cdot)$ is assumed to include the terms that depend on d_{\min} . This expression can be substituted in (7.5) to obtain

$$P_e[\mathcal{E}] \leq \sum_{N=1}^{\infty} \sum_{s' \in S'_N} \sum_{\hat{s}' \in S'_N(s')} \sum_{x \in \mathcal{X}_N(s')} P[x] \sum_{\substack{\hat{x} \in \mathcal{X}_N(\hat{s}') \\ \hat{x} \neq x}} \prod_{n=0}^{N-1} \sum_{z_n \in \mathcal{Z}} e^{q(x_n, \hat{x}_n, z_n)} P[z_n] \quad (7.35)$$

Furthermore, using the substitution $P[x] = P[s'] \prod_{n=0}^{N-1} P[x_n | s'_n]$, we get

$$P_e[\mathcal{E}] \leq \sum_{N=1}^{\infty} \sum_{s' \in S'_N} \sum_{\hat{s}' \in S'_N(s')} P[s'] \sum_{x \in \mathcal{X}_N(s')} \prod_{n=0}^{N-1} P[x_n | s'_n] \cdot \sum_{\substack{\hat{x} \in \mathcal{X}_N(\hat{s}') \\ \hat{x} \neq x}} \prod_{n=0}^{N-1} \sum_{z_n \in \mathcal{Z}} e^{q(x_n, \hat{x}_n, z_n)} P[z_n] \quad (7.36)$$

Let us assume that the symbols in the source sequence are mutually independent and have equiprobable values. Then, for an Ungerboeck code,

$$P[s'] = \begin{cases} (1/2^\nu) (1/2^{\tilde{m}})^N & \text{if } s' \in S'_N \\ 0 & \text{otherwise} \end{cases} \quad (7.37)$$

and

$$P[x_n | s'_n] = \begin{cases} 1/2^{m-\tilde{m}} & \text{if } x_n \in \mathcal{X}(s'_n) \\ 0 & \text{otherwise} \end{cases} \quad (7.38)$$

where $1/2^\nu$ is the probability of a given initial state, $1/2^{\tilde{m}}$ is the probability of transition to an allowable next state, and $2^{m-\tilde{m}}$ is the number of parallel transitions per state transition (all equiprobable). The union bound can now be expressed as

$$P_e[\mathcal{E}] \leq \frac{1}{2^\nu} \sum_{N=1}^1 \sum_{s' \in S'_N} \sum_{\hat{s}' \in S'_N(s')} \prod_{n=0}^{N-1} \frac{1}{2^{\tilde{m}}} \sum_{x_n \in \mathcal{X}(s'_n)} \sum_{\substack{\hat{x}_n \in \mathcal{X}(s'_n) \\ \hat{x}_n \neq x_n}} \sum_{z_n \in \mathcal{Z}} e^{q(x_n, \hat{x}_n, z_n)} P[z_n] \\ + \frac{1}{2^\nu} \sum_{N=2}^{\infty} \sum_{s' \in S'_N} \sum_{\hat{s}' \in S'_N(s')} \prod_{n=0}^{N-1} \frac{1}{2^{\tilde{m}}} \sum_{x_n \in \mathcal{X}(s'_n)} \sum_{\substack{\hat{x}_n \in \mathcal{X}(s'_n) \\ \hat{x}_n \neq x_n}} \sum_{z_n \in \mathcal{Z}} e^{q(x_n, \hat{x}_n, z_n)} P[z_n] \quad (7.39)$$

To evaluate this bound numerically, it is useful to write it in matrix notation [Biglieri, 1984]

$$P_e[\mathcal{E}] \leq T_e(\sigma_\eta^2, \rho) = \frac{1}{2^\nu} \mathbf{1}_G^t \mathbf{T}_{GG} \mathbf{1}_G + \frac{1}{2^\nu} \mathbf{1}_G^t \mathbf{T}_{GB} \left\{ \sum_{n=0}^{\infty} \mathbf{T}_{BB}^n \right\} \mathbf{T}_{BG} \mathbf{1}_G \quad (7.40)$$

where ρ is replaced by λ if the Chernoff bound is used rather than the Viterbi bound. The elements of the matrix \mathbf{T} are

$$\mathbf{T}_{(s_n, \hat{s}_n)(s_{n+1}, \hat{s}_{n+1})} = \begin{cases} \frac{1}{2^{\tilde{m}}} \sum_{x_n \in \mathcal{X}(s'_n)} \sum_{\substack{\hat{x}_n \in \mathcal{X}(s'_n) \\ \hat{x}_n \neq x_n}} \sum_{z_n \in \mathcal{Z}} e^{q(x_n, \hat{x}_n, z_n)} P[z_n] & \text{if } s'_n, \hat{s}'_n \in S' \\ 0 & \text{if } s'_n, \hat{s}'_n \notin S' \end{cases} \quad (7.41)$$

Subscripts G and B denote 'good' ($s_n = \hat{s}_n$) and 'bad' ($s_n \neq \hat{s}_n$) subsets of state pairs so that, for example, \mathbf{T}_{GG} is a submatrix of \mathbf{T} containing the contribution of error events

between parallel transitions to the error-event probability. The matrix $\mathbf{1}_G^t$ is a row vector with all elements unity, which sums over all possible initial 'good' state pairs, and the matrix $\mathbf{1}_G$ is a column vector, which sums over all possible final 'good' states. Provided $\mathbf{I}_{BB} - \mathbf{T}_{BB}$ is non-singular, where \mathbf{I}_{BB} is an identity matrix, the infinite matrix series in (7.40) can be written as

$$\left\{ \sum_{n=0}^{\infty} \mathbf{T}_{BB}^n \right\} = [\mathbf{I}_{BB} - \mathbf{T}_{BB}]^{-1} \quad (7.42)$$

Thus, (7.40) can be expressed as

$$T_e(\sigma_\eta^2, \rho) = \frac{1}{2^\nu} \mathbf{1}_G^t \left\{ \mathbf{T}_{GG} + \mathbf{T}_{GB} [\mathbf{I}_{BB} - \mathbf{T}_{BB}]^{-1} \mathbf{T}_{BG} \right\} \mathbf{1}_G \quad (7.43)$$

which is a generalized transfer function of the trellis code based on the $2^{2\nu}$ error states. This is usually denoted $T_e(D)$, but the parameter $D = \exp(-1/8\sigma_\eta^2)$ has been omitted because it is incompatible with the form of the Chernoff bound.

If a bound on the bit or symbol error probability is required, then the generalized transfer function must be modified to include a factor J

$$\mathbf{T}(J)_{(s_n, \hat{s}_n)(s_{n+1}, \hat{s}_{n+1})} = \frac{1}{2^m} \sum_{x_n \in \mathcal{X}(s'_n)} \sum_{\hat{x}_n \in \mathcal{X}(\hat{s}'_n)} \sum_{z_n \in \mathcal{Z}} J^k e^{q(x_n, \hat{x}_n, z_n)} P[z_n] \quad (7.44)$$

where for symbol error probability, the exponent of J is

$$k = \begin{cases} 0 & \text{if } x_n = \hat{x}_n \\ 1 & \text{if } x_n \neq \hat{x}_n \end{cases} \quad (7.45)$$

The bound on symbol error probability can now be computed from

$$P_s[\mathcal{E}] \leq T_s(\sigma_\eta^2, \rho) = \left. \frac{\partial T_s(\sigma_\eta, \rho, J)}{\partial J} \right|_{J=1} \quad (7.46)$$

where

$$\begin{aligned} \frac{\partial T_s(\sigma_\eta^2, \rho, J)}{\partial J} &= \frac{1}{2^\nu} \mathbf{1}_G^t \left\{ \mathbf{T}'_{GG} + \mathbf{T}_{GB} [\mathbf{I} - \mathbf{T}_{BB}]^{-1} \mathbf{T}'_{BB} [\mathbf{I} - \mathbf{T}_{BB}]^{-1} \mathbf{T}_{BG} \right. \\ &\quad \left. + \mathbf{T}_{GB} [\mathbf{I} - \mathbf{T}_{BB}]^{-1} \mathbf{T}'_{BG} + \mathbf{T}'_{GB} [\mathbf{I} - \mathbf{T}_{BB}]^{-1} \mathbf{T}_{BG} \right\} \mathbf{1}_G \end{aligned} \quad (7.47)$$

and

$$\mathbf{T}' \equiv \frac{\partial \mathbf{T}}{\partial J} \quad (7.48)$$

The exponential terms averaged over the ISI in (7.41) can be computed by averaging e^q over the pdf of the ISI. For small ratios of ISI variance to noise variance, it may be faster to express e^q as a Taylor series so that the average exponential is an infinite series involving moments of the ISI. The average exponential terms can then be approximated by truncating the series and computing the required moments using the techniques described by Cariolaro and Pupolin [1975]. But when the ratio of ISI variance to noise variance is large (this is one of the cases we are interested in), the series may oscillate [Ho and Yeh, 1970] so that a large number of terms may be necessary to compute a good approximation to the average exponential. Averaging over the ISI pdf is guaranteed to give a good approximation, provided a good approximation to the ISI pdf is available. Since we have a technique to obtain good approximations of ISI pdf's (see Chapter 6), we average over the ISI pdf.

In practice, $[\mathbf{I}_{BB} - \mathbf{T}_{BB}]^{-1}$ is evaluated by truncating the series expansion in (7.42). Also, the exponential terms in all the matrices are thresholded to include only the significant terms because they rapidly approach zero as the distance between state transitions increases (all terms of the order of $P[\mathcal{E}]$ should be included). To further speed up the computation, sparse matrix techniques [Pissanetzky, 1984] can be applied.

The evaluation of (7.43) or (7.47) generally requires an analysis of all possible pairs of code states (*pairwise states*), which involves the manipulation of $(2^{2\nu} - 2^\nu) \times (2^{2\nu} - 2^\nu)$ matrices and is limited in practice to codes with $2^\nu \leq 64$ states. For codes with the *uniform error property* [Benedetto *et al.*, 1988; Biglieri and McLane, 1989] and $h_0 = 1$, these matrices can be reduced to a maximum size of $(2^\nu - 1) \times (2^\nu - 1)$ by considering a modified transfer function based on the 2^ν code states [Zehavi and Wolf, 1987; Benedetto *et al.*, 1988].

7.5 Numerical Evaluation of the Minimum Distance

An asymptotic upper bound on the minimum distance is given by Divsalar [1978]

$$d_{min}^2(\rho) \leq \lim_{\sigma_\eta^2 \rightarrow 0} 16\sigma_\eta^2 \ln \left(\frac{T_e(2\sigma_\eta^2, \rho)}{T_e(\sigma_\eta^2, \rho)} \right) \quad (7.49)$$

An asymptotic lower bound on the minimum distance is given by Biglieri [1984]

$$d_{min}^2(\rho) \geq \lim_{\sigma_\eta^2 \rightarrow 0} -8\sigma_\eta^2 \ln \left(T_e(\sigma_\eta^2, \rho) \right) \quad (7.50)$$

but is slower to converge than the upper bound.

Small values of $d_{min}^2(\rho)$ are difficult to evaluate accurately using pairwise error probabilities obtained by averaging over the ISI pdf. The ISI values that govern $d_{min}^2(\rho)$ are those of greatest magnitude. Accurate values of $d_{min}^2(\rho)$ can be computed using pairwise error probabilities obtained by averaging over an *extreme ISI pdf*. The extreme ISI pdf consists only of ISI values with greatest magnitude ($\max\{\text{Re}[z_n]\}$). Because the ISI pdf's are usually circularly symmetric, an extreme ISI pdf consisting of a circle of eight impulses is sufficient to accurately compute $d_{min}^2(\rho)$.

7.6 Lower Bound on Error-Event Probability

An obvious lower bound on $P_e[\mathcal{E}]$ is the asymptotic lower bound of the Ungerboeck code in the absence of ISI (described in Section 3.4)

$$P_e[\mathcal{E}] \geq N(d_{free}) Q \left(\frac{d_{free}}{2\sigma_\eta} \right) \quad (7.51)$$

where d_{free} is the free distance of the code and $N(d_{free})$ is the average number of paths at distance d_{free} from a given reference path. This provides a tight bound for low levels of ISI, but is loose for severe ISI.

Forney gives tight lower bounds on error-event probability for the Viterbi algorithm

$$P_e[\mathcal{E}] \geq K_0 Q \left(\frac{d_{min}}{2\sigma_\eta} \right) \quad (7.52)$$

$$(7.53)$$

However, the computation of the K_0 parameter is not straightforward for our problem and we do not use these bounds.

7.7 Examples of Bounds on Error Probability

We now examine bounds on the error probability, which have been computed using the techniques just described, for two Ungerboeck coded systems. The first system is a $\nu = 2$ coded 16-QAM system that uses a rate $1/2$, 4 state convolutional encoder and 2 uncoded bits to map onto a 16-QAM signal constellation. This is a linear code, and is analyzed using a modified transfer function based on the code states. The second system is a $\nu = 3$ coded 512-CR system that uses a rate $2/3$, 8 state nonlinear convolutional encoder and 6 uncoded bits to map onto a 512-CR signal constellation. This nonlinear code was designed to satisfy Wei's conditions for 90° rotational invariance, and is analyzed using the pairwise state analysis. The minimum distance between signal points in the 16-QAM and 512-CR signal constellations is $d = 2$.

The systems in the examples incorporate raised-cosine Nyquist filtering with a roll-off factor of 0.3 and a Nyquist bandwidth of 17.5 MHz. The time-dispersive channel was modelled using the Rummmler model (see Section 4.2.1) with a delay $\tau = 6.3$ ns between the two paths. Because TCM cannot alone combat the ISI introduced by typically encountered fade parameters, we consider the residual ISI after equalization with a five-tap *synchronously-spaced equalizer*. The tap weights of this equalizer are optimized using the *minimum mean-square error* criterion (assuming high SNR). Noise enhancement due to the equalizer is accounted for by an additional fixed noise source.

The equalized channel impulse responses we consider all have $h_0 \approx 1$; therefore, the multiplicative interference is negligible. Table 7.1 lists the peak ISI values $\max\{\text{Re}[z_n]\} = \max\{\text{Im}[z_n]\}$, the optimum value of ρ at high SNR, and the ISI-degraded minimum distance values $d_{\min}(\rho)$, computed using (7.49), for the TCM schemes on various channels. Notice that the optimum value of ρ approaches zero as $\max\{\text{Re}[z_n]\}$ increases. However, because ρ is optimized over all conditional pairwise error probabilities, it is greater than zero for some values of $\max\{\text{Re}[z_n]\}$ greater than unity. The ISI pdf's have been computed using the central binning described in Chapter 6, with a bin width that ensures good approximations to the exact ISI pdf's. Worst case binning could also be used, but it would lead to a looser union bound.

The tightness of the upper bounds is determined by comparing them to the results of Monte Carlo simulations (see Chapter 5) and the lower bound. The systems have

System	B (dB)	f_o (MHz)	$\max\{\text{Re}[z_n]\}$	ρ	$d_{\min}(\rho)$
coded 16-QAM	—	—	0.00	0.50	4.0
	10	4	0.23	0.47	3.7
	15	4	0.41	0.44	3.3
	20	4	1.04	0.20	1.8
	25	4	2.16	0.00	0.0
coded 512-CR	—	—	0.00	0.50	4.5
	5	4	0.31	0.41	2.8
	10	0	0.50	0.40	2.8
	9.4	4	0.96	0.15	0.9
	16	0	1.16	0.10	0.0

Table 7.1 Peak ISI introduced by the channel and the ISI-degraded minimum distance of the TCM.

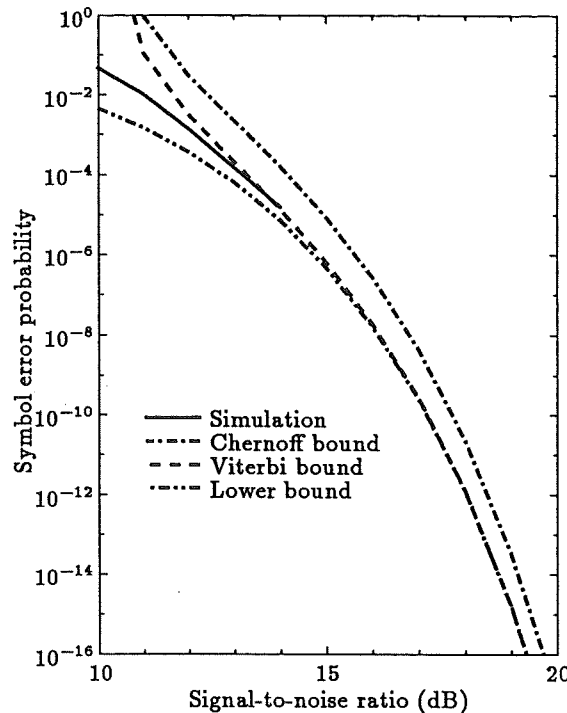


Figure 7.2 Symbol error probability bounds for coded 16-QAM on an AWGN channel.

been simulated both with and without interleaving to determine when the independence assumptions in Section 7.2 are valid; an interleaving depth of 20 symbols was used. With the resources available, simulations can only accurately estimate symbol error probabilities greater than 10^{-5} , but nonetheless the simulations provide an indication of how tight the bounds are. The simulations use a finite decoding depth in the Viterbi decoder of 12 symbols, whereas the bounds assume an infinite decoding depth.

Both optimized and unoptimized bounds are shown. The advantage of the unoptimized bounds is that they require less computation, while the advantage of the optimized bounds is that they are tighter than the unoptimized bounds for severe residual ISI. The Viterbi bound has been computed using either the computed ISI pdf (Viterbi-C) or assuming the ISI has a Gaussian pdf (Viterbi-G). The advantage of using a Gaussian ISI pdf is that the ISI is only specified by a mean and a variance, and the bound need not be optimized.

Bounds on the performance of the coded 16-QAM system as a function of SNR on an AWGN channel are shown in Figure 7.2. The lower bound is computed from (7.51) and is very tight relative to the simulation and the Viterbi bound (this same lower bound is used for the coded 16-QAM system with residual ISI). The Viterbi bound is also very tight relative to the simulation, for $P_s[\mathcal{E}] \leq 10^{-3}$. The Chernoff bound is about 0.5 dB looser than the Viterbi bound.

Figures 7.3a and b show bounds for the coded 16-QAM system for two residual ISI distributions with $\max\{\text{Re}[z_n]\} \leq 0.5$. The ISI pdf's can be found in Chapter 6. The lower bounds are looser than for the AWGN channel. Although the upper bounds are also looser than for an AWGN channel, they are still usefully tight for systems with interleaving. For these channel conditions, there is negligible difference between the unoptimized and optimized bounds. The upper bounds in Figure 7.3a apply to both interleaved and uninterleaved systems because interleaving provides negligible improvements in performance, as shown by the simulated curves. However, in Figure 7.3b the bounds only strictly apply to the interleaved system because interleaving provides a gain of 0.5 dB in noise margin.

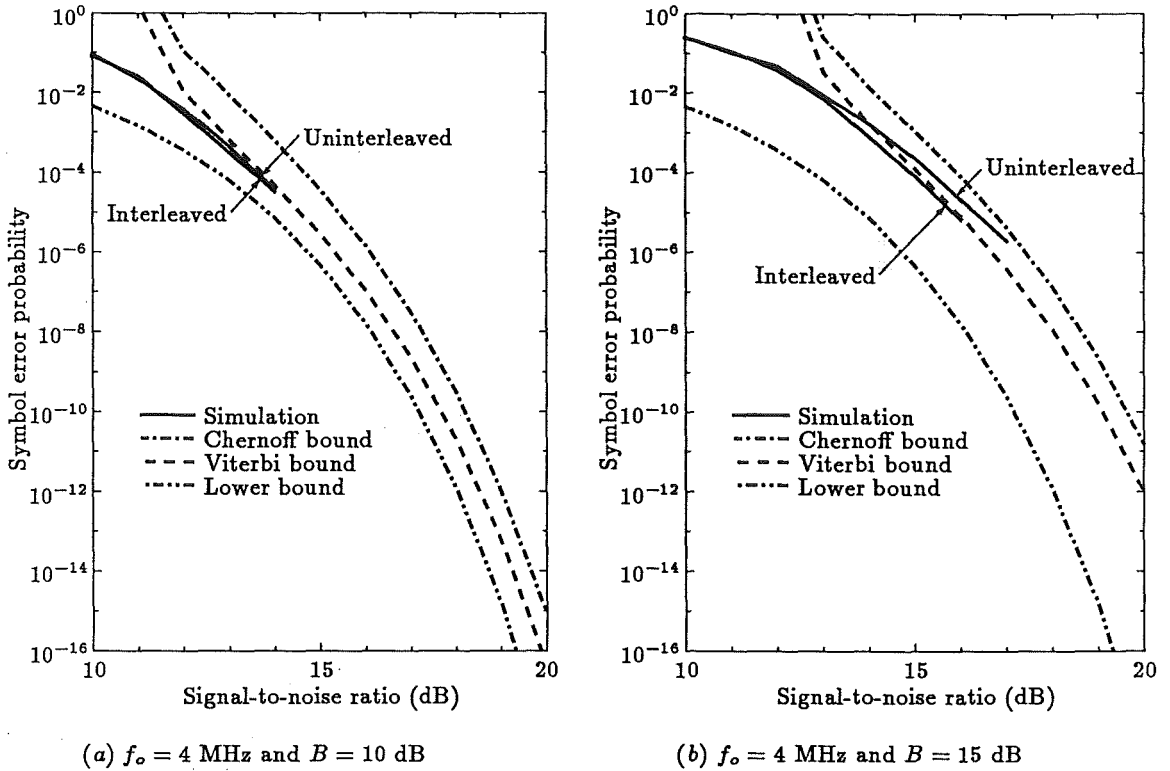


Figure 7.3 Symbol error probability bounds for coded 16-QAM on residual ISI channels with $\max\{\text{Re}[z_n]\} \leq 0.5$.

Applying the Viterbi bound using a Gaussian ISI approximation yields a bound that is indistinguishable from the Viterbi bound using the computed ISI pdf. This suggests that, for low levels of ISI, the ISI can be well approximated as a Gaussian random variable. In practice this means that an upper bound can be computed for low levels of ISI without having to compute an ISI pdf or having to optimize the bound.

Figures 7.4a and b show bounds for coded 16-QAM for two residual ISI distributions with $\max\{\text{Re}[z_n]\} > 0.5$. The optimized bounds in Figure 7.4a are tight; however, the unoptimized Viterbi bound is quite loose due to the failure of the bound when $\max\{\text{Re}[z_n]\} > 0.5$, as discussed in Section 7.3.1. The Viterbi bound with Gaussian ISI is slightly tighter than with the computed ISI for $S/N < 20$ dB, but becomes much looser for $S/N > 20$ dB. Thus, in this case, the ISI is not well approximated by a Gaussian random variable. The unoptimized Chernoff bound is becoming loose at an SNR of 22 dB, while the unoptimized Viterbi bound becomes loose for $S/N > 20$ dB. The convex (U) shape of this bound is due to the bound on a/\sqrt{b} becoming loose for many conditional pairwise error probabilities as the SNR increases. The interleaving provides a gain of about 1 dB in noise margin.

The residual ISI for Figure 7.4b was chosen to produce an error floor at an error probability that could be simulated. For this level of ISI, $d_{\min} = 0$. It was not possible to compute a meaningful Viterbi bound, either optimized or unoptimized, using a computed ISI pdf because the bound on a/\sqrt{b} was too loose for too many conditional pairwise error probabilities. However, it was possible to calculate the Viterbi bound using a Gaussian ISI approximation, and this is somewhat tighter than the Chernoff bounds. Thus, even though the ISI is not well approximated by a Gaussian random variable, a Viterbi-G bound on the error probability floor is just as tight as a Chernoff bound. The optimized Chernoff bound is one to two orders of magnitude different to the uninterleaved simulation, and therefore

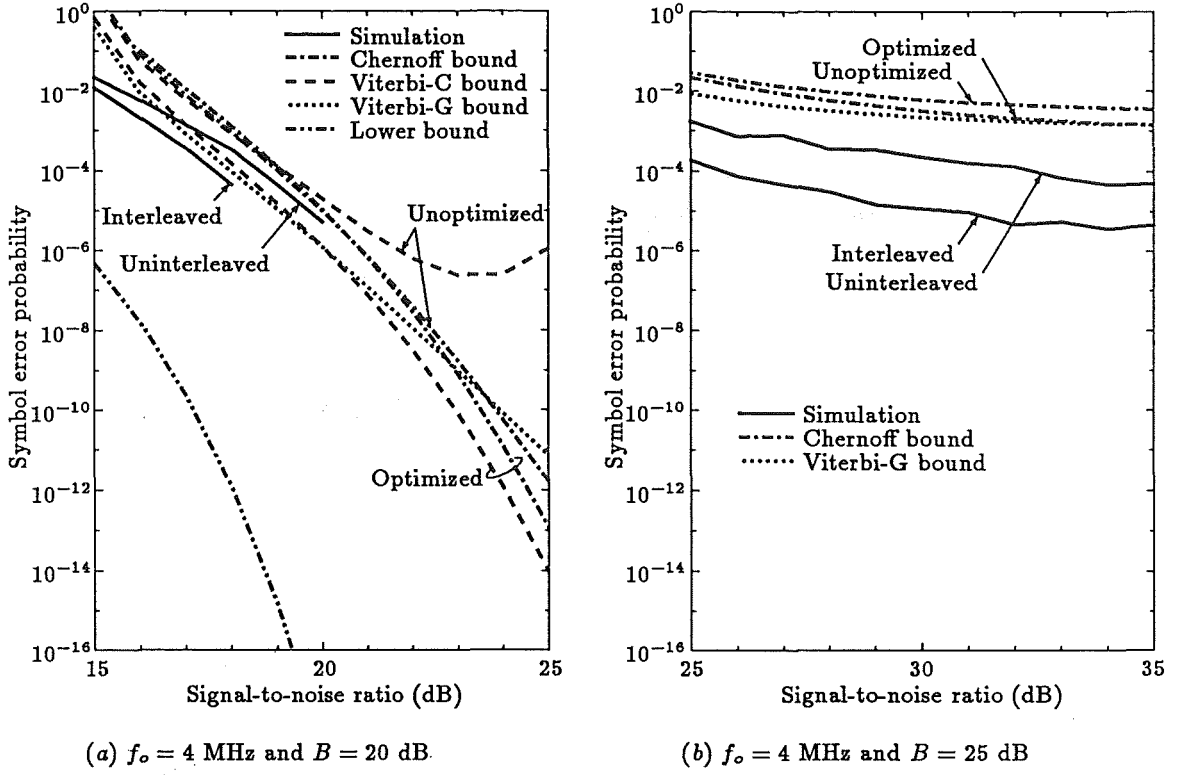


Figure 7.4 Symbol error probability bounds for coded 16-QAM on residual ISI channels with $\max\{\text{Re}[z_n]\} > 0.5$.

is quite loose. The unoptimized Chernoff bound is slightly looser again. Interleaving provides an order of magnitude gain in error probability compared to no interleaving. If a tighter bound on the error floor could be obtained, then it could be used to compute tight bounds on system signatures at low error probabilities (e.g. 10^{-6}), which cannot easily be simulated.

We now study the bounds for the coded 512-CR system. Bounds on the performance of the coded 512-CR system on an AWGN channel are shown in Figure 7.5. The lower bound is very tight relative to the simulation. As for the coded 16-QAM system, this lower bound is also used for the coded 512-CR system with residual ISI. The same observations as for the coded 16-QAM on an AWGN channel apply to coded 512-CR. Further to these observations we notice that the simulation curve crosses the Viterbi bound. It is suspected that this is due to the uncertainty associated with the points in the simulated curve. Although the uncertainty is difficult to compute analytically, the results in Chapter 5 suggest it could be up to 5%.

The bounds for coded 512-CR for two residual ISI distributions with $\max\{\text{Re}[z_n]\} \leq 0.5$ are shown in Figures 7.6a and b. The interleaving has little effect with either of these channel conditions. This is because, to yield the same peak ISI as with a coded 16-QAM system, the fades are significantly less severe. Also the coded 512-CR system uses a more powerful code than the coded 16-QAM system.

The bounds for two residual ISI distributions with $\max\{\text{Re}[z_n]\} > 0.5$ are shown in Figures 7.7a and b. The performance of the coded 512-CR system with these residual ISI conditions was studied by simulation in Chapter 5. As for the coded 16-QAM examples, we find that the optimized bounds are usefully tight, the unoptimized Chernoff bounds become looser at high SNRs, and the unoptimized Viterbi bounds are of no practical use. Interleaving has a negligible effect in Figure 7.7a, but provides about 0.5 dB improvement

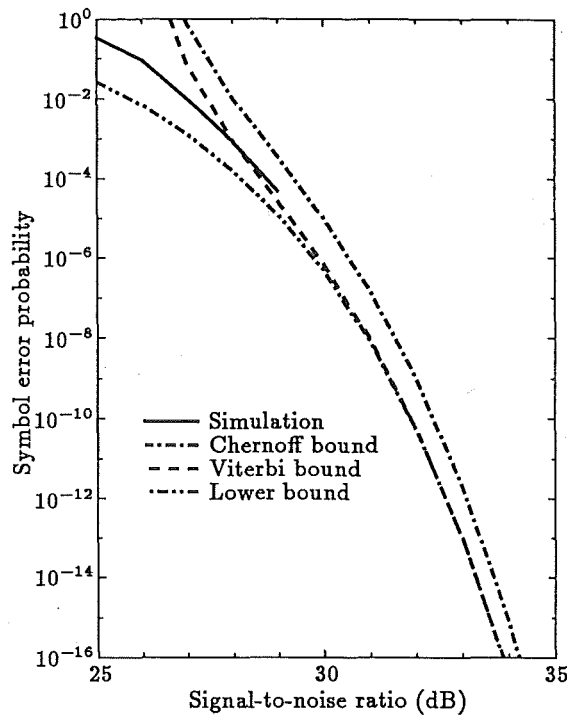
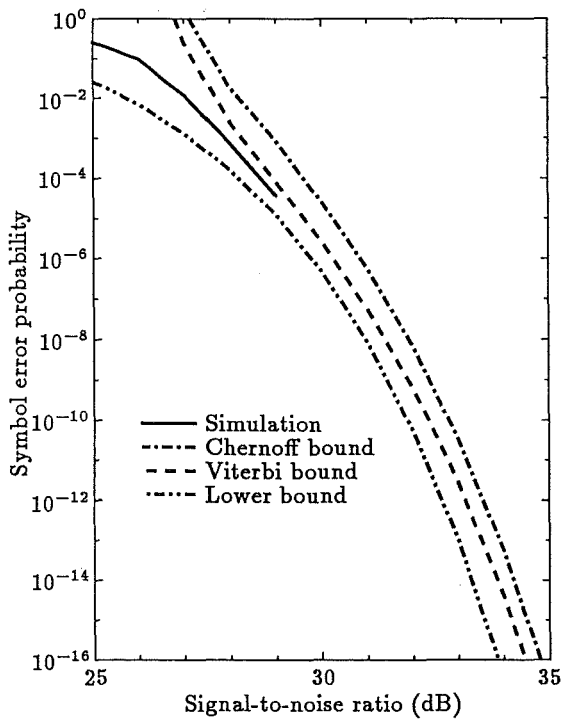
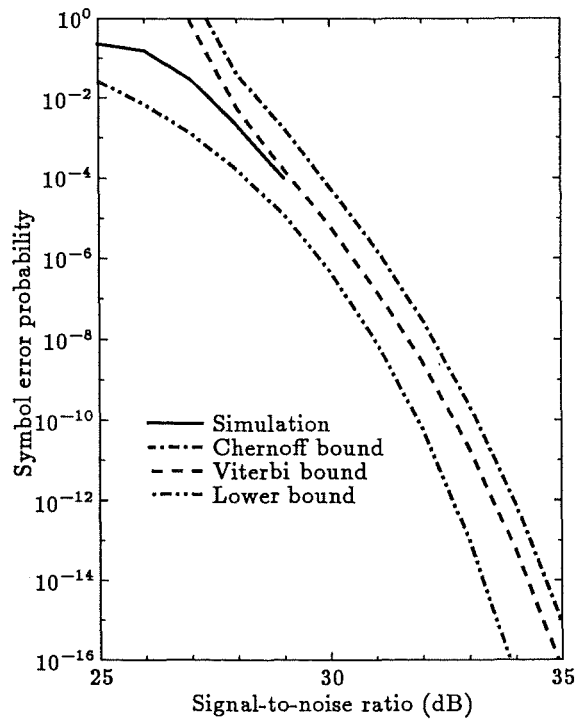


Figure 7.5 Symbol error probability bounds for coded 512-CR on an AWGN channel.



(a) $f_o = 4$ MHz and $B = 5$ dB



(b) $f_o = 0$ MHz and $B = 10$ dB

Figure 7.6 Symbol error probability bounds for coded 512-CR on residual ISI channels with $\max\{\text{Re}[z_n]\} \leq 0.5$.

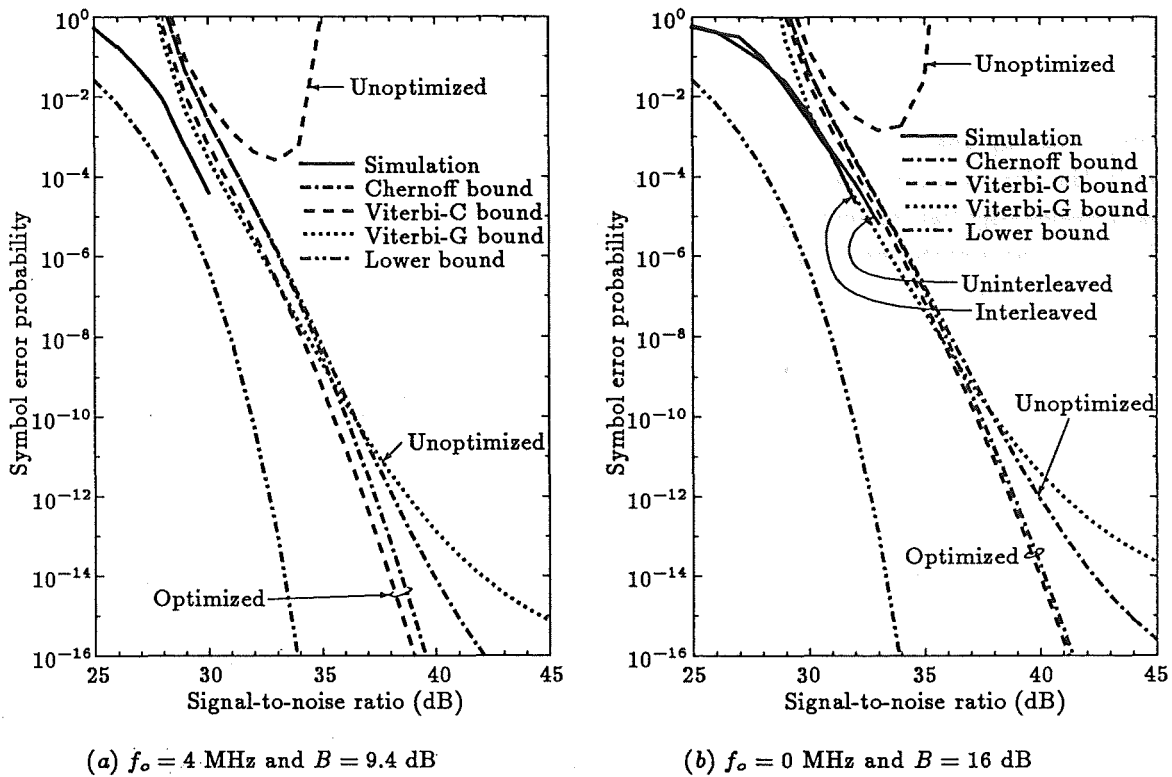


Figure 7.7 Symbol error probability bounds for coded 512-CR on residual ISI channels with $\max\{\text{Re}[z_n]\} > 0.5$.

in Figure 7.7b.

It should be mentioned that error probability curves measured for real radio systems usually show an error floor at $P_s[\mathcal{E}] = 10^{-14}$ due to irreducible sources of noise or interference. Such effects are due to modem imperfections and have not taken into account within the analysis in this chapter.

Analyzing the nonlinear Wei code using the code states gives virtually identical bounds to using pairwise states. This suggests that there is a larger class of codes that can be analyzed (or approximately analyzed) using the Zehavi and Wolf technique to avoid having to analyze pairwise states.

7.8 Conclusion

Upper bounds have been formulated and computed for the error probability of TCM on time-dispersive channels. The availability of accurate approximations to ISI pdf's has been used in the formulation. Examples have shown the bounds to be sufficiently tight over a wide range of time-dispersive channel conditions for practical applications. For severe levels of ISI, the bounds must be optimized over a bounding parameter to yield tight bounds; however, for low levels of ISI, a fixed bounding parameter can be used to yield tight bounds. A further simplification that can be used for low levels of ISI is to assume that the ISI has a Gaussian pdf, to avoid having to compute an approximate pdf.

The lower bound on the error probability of TCM on an AWGN channel has been used as a lower bound for TCM on a residual ISI channel. This bound is tight for low levels of ISI, but loose for severe ISI. Tighter lower bounds are required, so the tightness of the upper bounds can be verified at lower error probabilities than can be verified by

simulation. If tighter lower or upper bounds on error probability floors could be obtained, then system signatures could be estimated without simulation.

We have observed that interleaving has little effect on the performance of a TCM system experiencing shallow fades. As a consequence of this, the independence assumptions made in the derivation of the union bound are valid for shallow fades. However, for deep fades, interleaving can provide up to 1 dB of gain, and the union bound is not strictly applicable to an uninterleaved system.

The techniques in this chapter could be extended to analyze the effects of adjacent channel interference and co-channel interference (with some appropriate assumptions) on TCM, particularly in the cellular mobile radio environment.

Appendix 7A ISI Limit for Viterbi Bound

The Viterbi bound on the conditional pairwise error probability becomes loose when the ISI exceeds a limit. This limit can be approximately calculated for a trellis-coded QAM constellation with minimum distance d between signal points by considering the condition that determines the limit

$$-\mu_{\Delta M} = \rho \sigma_{\Delta M}^2 / 4\sigma_\eta^2 \quad (7.54)$$

The limit on the peak ISI depends on the pairwise error probability under consideration. We wish to derive the minimum value of the limit on the peak ISI for all pairwise errors. To do this we use the fact that when (7.54) is satisfied, the minimum distance $d_{min}(\rho) = 0$. The error event that will have this ISI-degraded minimum distance is the error event with minimum average squared Euclidean distance per branch. For the trellis-coded QAM constellations we consider, the minimum average squared distance per branch for an error event is $d^2/2$. The error event that has this distance is a very long event with the series of branch distances: $2d^2 + 0 + d^2 + 0 + d^2 + 0 + d^2 + \dots + 2d^2$.

The condition in (7.54) can be expanded and rearranged to obtain

$$(1 - \rho) \sum_{n=0}^{N-1} |x_n - \hat{x}_n|^2 = -2 \sum_{n=0}^{N-1} \text{Re} [(x_n - \hat{x}_n)^* z_n] \quad (7.55)$$

where we assume $h_0 = 1$. Knowing that the branch distances in the error event of interest alternate between 0 and d^2 , and assuming the $2d^2$ contributions of the first and last branches are negligible over length of the error event, we can use (7.55) to obtain

$$(1 - \rho) \frac{N}{2} d^2 = 2 \max \left\{ \frac{N}{2} \text{Re} [dz_n] \right\} \quad (7.56)$$

which can be solved to get a limit on the ISI

$$z_{max} = \max \{ \text{Re} [z_n] \} = (1 - \rho) d / 2 \quad (7.57)$$

This is the condition given in (7.26). An identical limit can be determined for $\max \{ \text{Im} [z_n] \}$.

Chapter 8

Conclusions

This thesis has achieved two major goals. First, the performance of trellis-coded modulation (TCM) on the time-dispersive digital microwave radio (DMR) channel has been evaluated. Second, tight analytical bounds on error probability that can be used to study the performance of TCM on general time-dispersive channels at low error probabilities, for which simulation is impractical, have been developed. In Chapter 5, Monte Carlo *simulation* techniques were used to study the performance of TCM on the DMR channel. Chapters 6 and 7 presented *analytical* techniques to study the performance of TCM on general time-dispersive channels. Detailed conclusions are given at the end of each chapter. This chapter provides more general conclusions and suggestions for further research.

The simulation study in Chapter 5 has clearly shown the significant improvements in performance that TCM can offer for DMR systems with equalization. Symbol error rate (SER) versus signal-to-noise ratio (SNR) curves, and computed outage probabilities were used to study the improvements in performance. Significant coding and SER gains were obtained with TCM over a range of channel conditions. Improvements in the severely errored seconds component of outage ($\text{BER} = 10^{-3}$) were not always obtained with TCM, but the improvements observed in outage for $\text{BER} = 10^{-4}$ suggest that the degraded minutes ($\text{BER} = 10^{-6}$) component of outage would be significantly improved with TCM. The results of this study are also applicable to TCM on other time-dispersive channels.

An interesting additional observation from the simulation study was the poor performance of the MMSE-DFE with the large signal constellations. This was first noted by Carlisle *et al.* [1989], and its cause was traced to attenuation of the cursor in the equalized impulse response.

In the past, when the probability density function (pdf) of the intersymbol interference (ISI) could not be computed exactly, researchers have either assumed that the pdf conforms to a uniform or a Gaussian pdf, or they have computed an approximation. However, the uniform and Gaussian approximations are often not accurate, and prior to the work in this thesis, the existing algorithms for computing approximations were not applicable to trellis-coded data. Chapter 6 presented algorithms for computing an approximation to the pdf of the ISI for uncoded and trellis-coded systems. These algorithms are based on binning a number of partial pdf's that can then be convolved to get an approximation to the ISI pdf. Worst and best case binning techniques were used to obtain ISI pdf's that gave upper and lower bounds on the symbol error probability. The examples presented showed that the computed ISI pdf's are good approximations to the actual ISI pdf's. The dependence between symbols, introduced by Ungerboeck codes, was found to have a negligible effect on the ISI pdf, so the algorithm for uncoded modulation can also be used

to approximate the ISI pdf's of Ungerboeck codes.

Few analytical techniques are available for computing bounds on the error probability of TCM on time-dispersive channels. Prior to the work in this thesis, the existing techniques were based on explicitly analyzing the channel states in addition to the code states. Because of the extensive computation that is involved, such techniques restrict the size of the signal constellation and length of the channel impulse response that can be analyzed. Chapter 7 developed bounds, on the error probability of TCM on time-dispersive channels, that do not require the channel states to be analyzed explicitly. The upper bounds were union bounds, formulated to make use of the approximate ISI pdf. Examples showed that these upper bounds are tight for a wide range of channel conditions. They also showed that assuming the ISI has a Gaussian pdf is reasonable when the ISI is small compared to the noise, but can lead to very loose bounds when the ISI is severe. The lower bound was simply the performance of the TCM with AWGN alone; this bound is tight for low levels of ISI, but loose for severe ISI.

8.1 Suggestions for Further Research

A number of ideas for further research arise from this work. Prior to a practical application of TCM to DMR systems, the effects of modem imperfections, in conjunction with residual ISI and noise, on the performance of TCM should be studied. Some common modem imperfections are carrier recovery errors, timing jitter, and high power amplifier nonlinearities.

The receivers we have considered are suboptimum because they only consider the code states and ignore the channel states. It would be useful to study the performance gains that can be obtained by applying TCM to time-dispersive channels and using the optimum receiver originally described by Forney [1972]. Forney's receiver must be simplified in practice, by ignoring some of the channel states. Therefore, it would also be useful to look at practical reduced complexity receivers. In particular, the work of Eyuboğlu and Qureshi [1989] and Chevillat and Eleftheriou [1989] on reduced complexity receivers for TCM on time-dispersive channels could be studied in the context of DMR systems. However, implementing these reduced complexity receivers at DMR data rates is wishful thinking with current technology, because a Viterbi decoder that just operates on code states is only marginally practical.

Since we are interested in TCM schemes that can be implemented at DMR data rates, there is little point in looking at more complicated codes. However, multidimensional trellis codes [Wei, 1987] and multilevel trellis codes [Calderbank, 1989; Pottie and Taylor, 1989a] may offer benefits over Ungerboeck codes in terms of reduced complexity.

The negligible effect that the dependence between symbols, introduced by Ungerboeck codes, has on the ISI pdf seems to relate to the symmetries built into the trellis structure of Ungerboeck codes. It also seems to relate to the null effect that Ungerboeck codes have on the power spectrum of the uncoded data [Biglieri, 1984]. Some very interesting work could be conducted in this area. A starting point could be trellis codes that introduce spectral nulls [Calderbank *et al.*, 1988; Wolf and Ungerboeck, 1986]. These codes are expected to have some effect on the ISI pdf, which could be computed using the techniques described in Chapter 6. This work could also look at other trellis codes and even error-control codes in general.

Moments of a random variable are often easier to compute than the pdf of the random variable. Cariolaro and Pupolin [1975] have described a technique to compute the ISI

moments for data that is encoded according to the rules of a finite state machine (this covers TCM). The union bound described in Chapter 7 could easily be reformulated to make use of the ISI moments rather than ISI pdf's, using similar techniques to Ho and Yeh [1970]. The use of ISI moments rather than ISI pdf's would reduce the cost of computing the error probability bounds when the variance of the ISI is of the order of the variance of the noise. However, when the variance of the ISI is larger than that of the noise, a large number of ISI moments would be required for the bound to converge [Ho and Yeh, 1970]. In this case the use of ISI pdf's is better.

If tighter bounds on error probability floors could be obtained, then system signatures could be estimated without simulation. Outage probabilities could then be estimated using these signatures, as in Chapter 5. This would be particularly useful to estimate the level of degraded minutes for a system. Tighter lower bounds on error probability would also be useful, especially to verify the tightness of the upper bounds for lower error probabilities than can be estimated by simulation.

Further work is required to develop bounds that do not require assumptions of independence between the ISI samples. The bounds described by Oka and Biglieri [1989] are of this type, but they consider all channel and code states (bounds that ignore some of the channel states are also described, but these bounds are loose), and can only be computed for simple TCM schemes and short channel impulse responses. A hybrid technique that considers a small number of channel states, to model the significant dependence between ISI samples, and uses a pdf based technique to model the remainder of the ISI may be possible. If such bounds were available, then there would be no need to assume that the systems we analyze include interleaving.

Finally, it would be useful to investigate the possibility of designing TCM specifically to cope with channels that introduce both residual ISI and AWGN. The analytical bounds on error probability that have been presented provide some useful tools for this task. Divsalar and Simon [1988a; 1988b] have designed codes for multiplicative interference channels that perform better than if codes designed for AWGN channels were used. Eyuboğlu [1988] shows that at high signal-to-noise ratios, if a coded modulation scheme can approach channel capacity on an ISI-free channel, then under mild assumptions the same scheme can also approach capacity on a channel with ISI, provided the receiver uses ideal decision-feedback equalization followed by maximum-likelihood decoding. Such findings suggest that TCM designed for AWGN channels cannot be greatly improved, in a practical sense, for ISI channels, and that the codes we have studied are nearly as good as we can expect.

References

- Amitay, N. and Greenstein, L.J. (1984), 'Multipath outage performance of digital radio receivers using finite-tap adaptive equalizers', *IEEE Trans. Commun.*, Vol. COM-32, No. 5, May, pp. 597-608.
- Belfiore, C.A. and Park, J.H., Jr. (1979), 'Decision-feedback equalization', *Proc. IEEE*, Vol. 67, No. 8, August, pp. 1143-1156.
- Benedetto, S., Marsan, M.A., Albertengo, G. and Giachin, E. (1988), 'Combined coding and modulation: theory and applications', *IEEE Trans. Inf. Theory*, Vol. 34, No. 2, March, pp. 223-236.
- Biglieri, E. (1984), 'High-level modulation and coding for nonlinear satellite channels', *IEEE Trans. Commun.*, Vol. COM-32, No. 5, May, pp. 616-626.
- Biglieri, E. and McLane, P.J. (1989), 'Uniform distance and error probability properties of TCM schemes', in *ICC'89 Conf. Rec.*, IEEE, Boston, Mass., pp. 45.2.1-45.2.6.
- Calderbank, A.R. (1989), 'Multilevel codes and multistage decoding', *IEEE Trans. Commun.*, Vol. 37, No. 3, March, pp. 222-229.
- Calderbank, A.R. and Sloane, N.J.A. (1987), 'New trellis codes based on lattices and cosets', *IEEE Trans. Inf. Theory*, Vol. IT-33, No. 2, March, pp. 177-195.
- Calderbank, A.R., Lee, T. and Mazo, J.E. (1988), 'Baseband trellis codes with a spectral null at zero', *IEEE Trans. Inf. Theory*, Vol. 34, No. 3, May, pp. 425-434.
- Cariolaro, G.L. and Pupolin, S.G. (1975), 'Moments of correlated digital signals for error probability evaluation', *IEEE Trans. Inf. Theory*, Vol. IT-21, No. 5, September, pp. 558-568.
- Carlisle, C.J., Kennedy, W.K. and Shafi, M. (1989), 'Outage simulations for digital microwave radio systems with trellis-coded modulation', in *ICC'89 Conf. Rec.*, Boston, Mass., pp. 33.2.1-33.2.5.
- Carlisle, C.J., Shafi, M. and Kennedy, W.K. (1990a), 'Trellis-coded modulation on digital microwave radio systems—Simulations for multipath fading channels', accepted for publication in *IEEE Trans. Commun.*
- Carlisle, C.J., Taylor, D.P., Kennedy, W.K. and Shafi, M. (1990b), 'The probability density of intersymbol interference for trellis-coded modulation', submitted to *IEEE Trans. Commun.*
- Carlisle, C.J., Taylor, D.P., Kennedy, W.K. and Shafi, M. (1990c), 'Performance bounds for trellis-coded modulation on time-dispersive channels', submitted to *IEEE Trans. Commun.*
- CCITT (1988), *Recommendation G.821: Error Performance of an International Digital Connection Forming Part of an ISDN*, Geneva.
- Chamberlain, J.K., Clayton, F.M., Sari, H. and Vandamme, P. (1986), 'Receiver techniques for microwave digital radio', *IEEE Commun. Mag.*, Vol. 24, No. 11, November, pp. 43-54.
- Chen, S.X., Smith, P.J., Shafi, M. and Vere-Jones, D. (1990), 'Some improvements to conventional importance sampling techniques for coded systems using Viterbi decoding', *Electronic Letters*, Vol. 26, No. 12, June, pp. 802-804.

- Chevillat, P.R. and Eleftheriou, E. (1989), 'Decoding of trellis-encoded signals in the presence of intersymbol interference and noise', *IEEE Trans. Commun.*, Vol. 37, No. 7, July, pp. 669–676.
- Chouly, A. and Sari, H. (1988), 'Application of trellis coding to digital microwave radio', in *ICC'88 Conf. Rec.*, pp. 15.1.1–15.1.5.
- Despinic, M., Kirkland, W.K. and Taylor, D.P. (1989), 'On the performance of trellis-coded modulation in digital microwave radio', in *ICC'89 Conf. Rec.*, Boston, Mass., pp. 33.3.1–33.3.5.
- Divsalar, D. (1978), *Performance of Mismatched Receivers on Bandlimited Channels*, Ph.D. thesis, University of California, Los Angeles.
- Divsalar, D. and Simon, M.K. (1988a), 'The design of trellis coded MPSK for fading channels: Performance criteria', *IEEE Trans. Commun.*, Vol. 36, No. 9, September, pp. 1004–1012.
- Divsalar, D. and Simon, M.K. (1988b), 'The design of trellis coded MPSK for fading channels: Set partitioning for optimum code design', *IEEE Trans. Commun.*, Vol. 36, No. 9, September, pp. 1013–1021.
- Eyuboğlu, M.V. (1988), 'Detection of coded modulation signals on linear, severely distorted channels using decision-feedback noise prediction with interleaving', *IEEE Trans. Commun.*, Vol. 36, No. 4, April, pp. 401–409.
- Eyuboğlu, M.V. and Qureshi, S.U.H. (1988), 'Reduced-state sequence estimation with set partitioning and decision feedback', *IEEE Trans. Commun.*, Vol. 36, No. 1, January, pp. 13–20.
- Eyuboğlu, M.V. and Qureshi, S.U.H. (1989), 'Reduced-state sequence estimation for coded modulation on intersymbol interference channels', *IEEE J. Sel. Areas Commun.*, Vol. SAC-7, No. 6, August, pp. 989–995.
- Falconer, D.D. and Magee, F.R., Jr. (1973), 'Adaptive channel memory truncation for maximum likelihood sequence estimation', *Bell Syst. Tech. J.*, Vol. 52, No. 9, November, pp. 1541–1562.
- Forney, G.D., Jr. (1972), 'Maximum-likelihood sequence estimation of digital sequences in the presence of intersymbol interference', *IEEE Trans. Inf. Theory*, Vol. IT-18, No. 3, May, pp. 363–378.
- Forney, G.D., Jr. (1973), 'The Viterbi algorithm', *Proc. IEEE*, Vol. 61, No. 3, March, pp. 268–278.
- Forney, G.D., Jr. (1988a), 'Coset codes—Part I: Introduction and geometrical classification', *IEEE Trans. Inf. Theory*, Vol. 34, No. 5, September, pp. 1123–1151.
- Forney, G.D., Jr. (1988b), 'Coset codes—Part II: Binary lattices and related codes', *IEEE Trans. Inf. Theory*, Vol. 34, No. 5, September, pp. 1152–1187.
- Forney, G.D., Jr., Gallager, R.G., Lang, G.R., Longstaff, F.M. and Qureshi, S.U. (1984), 'Efficient modulation for band-limited channels', *IEEE J. Sel. Areas Commun.*, Vol. SAC-2, No. 5, September, pp. 632–646.
- Foschini, G.J. (1975), 'Performance bound for maximum-likelihood reception of digital data', *IEEE Trans. Inf. Theory*, Vol. IT-21, No. 1, January, pp. 47–50.
- Franks, L.E. (1980), 'Carrier and bit synchronization in data communication—A tutorial review', *IEEE Trans. Commun.*, Vol. COM-28, No. 8, August, pp. 1107–1121.
- Glave, F.E. (1972), 'An upper bound on the probability of error due to intersymbol interference for correlated digital signals', *IEEE Trans. Inf. Theory*, Vol. IT-18, No. 3, May, pp. 356–363.
- Greenstein, L.J. and Shafi, M. (1987), 'Outage calculation methods for microwave digital radio', *IEEE Commun. Mag.*, Vol. 25, No. 2, February, pp. 30–39.
- Hayes, J.F. (1975), 'The Viterbi algorithm applied to digital data transmission', *IEEE Commun. Mag.*, Vol. 13, No. 2, March, pp. 15–20.
- Haykin, S. (1986), *Adaptive Filter Theory*, Prentice-Hall, N.J.

- Heller, J.A. and Jacobs, I.M. (1971), 'Viterbi decoding for satellite and space communication', *IEEE Trans. Commun. Tech.*, Vol. COM-19, No. 5, October, pp. 835-848.
- Herro, M.A. and Nowack, J.M. (1988), 'Simulated Viterbi decoding using importance sampling', *Proc. IEE, Pt. F*, Vol. 135, No. 2, April, pp. 133-142.
- Hill, F.S., Jr. (1971), 'The computation of error probability for digital transmission', *Bell Syst. Tech. J.*, Vol. 50, No. 6, July-August, pp. 2055-2077.
- Hill, F.S., Jr. and Blanco, M.A. (1973), 'Random geometric series and intersymbol interference', *IEEE Trans. Inf. Theory*, Vol. IT-19, No. 3, May, pp. 326-335.
- Ho, E.Y. and Yeh, Y.S. (1970), 'A new approach for evaluating the error probability in the presence of intersymbol interference and additive Gaussian noise', *Bell Syst. Tech. J.*, Vol. 49, November, pp. 2249-2265.
- Huzii, A. and Sugiyama, H. (1970), 'Intersymbol interference of markov pulse trains', *Electronics and Communications in Japan*, Vol. 53-A, No. 3, pp. 21-30.
- Ivanek, F. (Ed.) (1989), *Terrestrial Digital Microwave Communications*, Artech House, Norwood, Mass.
- Jeruchim, M.C. (1984), 'Techniques for estimating the bit error rate in the simulation of digital communication systems', *IEEE J. Sel. Areas Commun.*, Vol. SAC-2, No. 1, January, pp. 153-170.
- Komaki, S., Nakamura, Y. and Kohiyama, K. (1990), 'Recent R & D on digital microwave radio relay systems', in *European Microwave Conf. Proc.*, pp. 107-119.
- Lin, S.H., Lee, T.C. and Gardina, M.F. (1988), 'Diversity protections for digital radio—Summary of ten-year experiments and studies', *IEEE Commun. Mag.*, Vol. 26, No. 2, February, pp. 51-64.
- Lucky, R.W. (1965), 'Automatic equalization for digital communications', *Bell Syst. Tech. J.*, Vol. 44, April, pp. 547-588.
- Lundgren, C.W. and Rummmler, W.D. (1979), 'Digital radio outage due to selective fading—Observation vs prediction from laboratory simulation', *Bell Syst. Tech. J.*, Vol. 58, No. 5, May-June, pp. 1073-1100.
- Magee, F.R., Jr. and Proakis, J.G. (1973), 'Adaptive maximum-likelihood sequence estimation for digital signaling in the presence of intersymbol interference', *IEEE Trans. Inf. Theory*, Vol. IT-19, No. 1, January, pp. 120-124.
- Matthews, J.W. (1973), 'Sharp error bounds for intersymbol interference', *IEEE Trans. Inf. Theory*, Vol. IT-19, No. 4, July, pp. 440-447.
- McKay, R.G. and Shafi, M. (1988), 'Multipath propagation measurements on an overwater path in New Zealand', *IEEE Trans. Commun.*, Vol. 36, No. 7, July, pp. 781-788.
- Metzger, K. (1987), 'On the probability density of intersymbol interference', *IEEE Trans. Commun.*, Vol. COM-35, No. 4, April, pp. 396-402.
- Moridi, S. and Sari, H. (1985), 'Analysis of four decision-feedback carrier recovery loops in the presence of intersymbol interference', *IEEE Trans. Commun.*, Vol. COM-33, No. 6, June, pp. 543-550.
- Noguchi, T., Daido, Y. and Nossek, J.A. (1986), 'Modulation techniques for microwave digital radio', *IEEE Commun. Mag.*, Vol. 24, No. 10, October, pp. 21-30.
- Nyquist, H. (1928), 'Certain topics in telegraph transmission theory', *Trans. AIEE*, Vol. 47, April, pp. 617-644.
- Oka, I. and Biglieri, E. (1989), 'Error probability bounds for trellis coded modulation over sequence dependent channels', *Transactions of the IEICE*, Vol. E 72, No. 4, April, pp. 375-383.

- Oppenheim, A.V. and Schaffer, R.W. (1975), *Digital Signal Processing*, Prentice-Hall, N.J.
- Pahlavan, K. and Holsinger, J.L. (1987), 'Expanded TCM for channels with multiplicative noise', in *ICC'87 Conf. Rec.*, IEEE, pp. 404-408.
- Pierce, J.R. (1980), *An Introduction to Information Theory—Symbols, Signals and Noise*, Dover, New York, 2nd ed.
- Pissanetzky, S. (1984), *Sparse Matrix Technology*, Academic Press.
- Pottie, G.J. and Taylor, D.P. (1987), 'An approach to Ungerboeck coding for rectangular signal sets', *IEEE Trans. Inf. Theory*, Vol. IT-33, No. 2, March, pp. 285-290.
- Pottie, G.J. and Taylor, D.P. (1989a), 'Multi-level codes based on partitioning', *IEEE Trans. Inf. Theory*, Vol. IT-35, No. 1, January, pp. 87-98.
- Pottie, G.J. and Taylor, D.P. (1989b), 'A comparison of reduced complexity decoding algorithms for trellis codes', *IEEE J. Sel. Areas Commun.*, Vol. 7, No. 9, December, pp. 1369-1380.
- Proakis, J.G. (1989), *Digital Communications*, McGraw-Hill, New York, 2nd ed.
- QUALCOMM (1990), '25 mbps QUALCOMM viterbi encoder', *IEEE Commun. Mag.*, Vol. 28, No. 3, March, p. 3. QUALCOMM incorporated: 10555 Sorrento Valley Rd., San Diego, Calif.
- Qureshi, S.U.H. (1985), 'Adaptive equalization', *Proc. IEEE*, Vol. 73, No. 9, September, pp. 1349-1387.
- Qureshi, S.U.H. and Newhall, E.E. (1973), 'An adaptive receiver for data transmission over time-dispersive channels', *IEEE Trans. Inf. Theory*, Vol. IT-19, No. 4, July, pp. 448-457.
- Rummler, W.D. (1979), 'A new selective fading model: Application to propagation data', *Bell Syst. Tech. J.*, Vol. 58, No. 5, May-June, pp. 1037-1071.
- Rummler, W.D., Coutts, R.P. and Liniger, M. (1986), 'Multipath fading channel models for microwave digital radio', *IEEE Commun. Mag.*, Vol. 24, No. 11, November, pp. 30-42.
- Saltzberg, B.R. (1968), 'Intersymbol interference error bounds with application to ideal bandlimited signaling', *IEEE Trans. Inf. Theory*, Vol. IT-14, No. 4, July, pp. 563-568.
- Schilling, D.L., Pickholtz, R.L. and Milstein, L.B. (1990), 'Spread spectrum goes commercial', *IEEE Spectrum*, Vol. 27, No. 8, August, pp. 40-45.
- Shannon, C.E. (1948a), 'A mathematical theory of communication', *Bell Syst. Tech. J.*, Vol. 27, No. 3, July, pp. 379-423. Parts I and II.
- Shannon, C.E. (1948b), 'A mathematical theory of communication', *Bell Syst. Tech. J.*, Vol. 27, October, pp. 623-656. Parts III, IV, and V.
- Shannon, C.E. (1949), 'Communication in the presence of noise', *Proc. IRE*, Vol. 37, January, pp. 10-21.
- Siller, C.A., Jr. (1984), 'Multipath propagation', *IEEE Commun. Mag.*, Vol. 22, No. 2, February, pp. 6-15.
- Sklar, B. (1988), *Digital Communications*, Prentice-Hall, Englewood Cliffs, N.J.
- Taylor, D.P. and Chan, H.C. (1981), 'A simulation study of two bandwidth-efficient modulation techniques', *IEEE Trans. Commun.*, Vol. COM-29, No. 3, March, pp. 267-275.
- Taylor, D.P. and Hartmann, P.R. (1986), 'Telecommunications by microwave digital radio', *IEEE Commun. Mag.*, Vol. 24, No. 8, August, pp. 11-16.
- Thapar, H.K. (1984), 'Real-time application of trellis coding to high-speed voiceband data transmission', *IEEE J. Sel. Areas Commun.*, Vol. SAC-2, No. 5, September, pp. 648-658.
- Townsend, A.A.R. (1988), *Digital Line-of-Sight Radio Links—A Handbook*, Prentice-Hall International (UK) Ltd.

- Ungerboeck, G. (1974), 'Adaptive maximum-likelihood receiver for carrier-modulated data-transmission systems', *IEEE Trans. Commun.*, Vol. COM-22, No. 5, May, pp. 624-636.
- Ungerboeck, G. (1982), 'Channel coding with multilevel/phase signals', *IEEE Trans. Inf. Theory*, Vol. IT-28, No. 1, January, pp. 55-67.
- Ungerboeck, G. (1987a), 'Trellis-coded modulation with redundant signal sets—Part II: State of the art', *IEEE Commun. Mag.*, Vol. 25, No. 2, February, pp. 12-21.
- Ungerboeck, G. (1987b), 'Trellis-coded modulation with redundant signal sets—Part I: Introduction', *IEEE Commun. Mag.*, Vol. 25, No. 2, February, pp. 5-11.
- Viterbi, A.J. (1967), 'Error bounds for convolutional codes and an asymptotically optimum decoding algorithm', *IEEE Trans. Inf. Theory*, Vol. IT-13, No. 2, April, pp. 260-269.
- Viterbi, A.J. (1971), 'Convolutional codes and their performance in communication systems', *IEEE Trans. Commun. Tech.*, Vol. COM-19, No. 5, October, pp. 751-772.
- Viterbi, A.J. and Omura, J.K. (1979), *Principles of Digital Communication and Coding*, McGraw-Hill, New York.
- Viterbi, A.J., Wolf, J.K., Zehavi, E. and Padovani, R. (1989), 'A pragmatic approach to trellis-coded modulation', *IEEE Commun. Mag.*, Vol. 27, No. 7, July, pp. 11-19.
- Wei, L.F. (1984a), 'Rotationally invariant convolutional channel coding with expanded signal space—Part I: 180°', *IEEE J. Sel. Areas Commun.*, Vol. SAC-2, No. 5, September, pp. 659-671.
- Wei, L.F. (1984b), 'Rotationally invariant convolutional channel coding with expanded signal space—Part II: Nonlinear codes', *IEEE J. Sel. Areas Commun.*, Vol. SAC-2, No. 5, September, pp. 672-686.
- Wei, L.F. (1987), 'Trellis-coded modulation with multidimensional constellations', *IEEE Trans. Inf. Theory*, Vol. IT-33, No. 4, July, pp. 483-501.
- Wolf, J.K. and Ungerboeck, G. (1986), 'Trellis coding for partial-response channels', *IEEE Trans. Commun.*, Vol. COM-34, No. 8, August, pp. 765-773.
- Wong, L.N. and McLane, P.J. (1988), 'Performance of trellis codes for a class of equalized ISI channels', *IEEE Trans. Commun.*, Vol. 36, No. 12, December, pp. 1330-1336.
- Wozencraft, J.M. and Jacobs, I.M. (1965), *Principles of Communication Engineering*, Wiley, New York.
- Yamamoto, H. (1987), 'Future trends in microwave digital radio', *IEEE Commun. Mag.*, Vol. 25, No. 2, February, pp. 40-52.
- Zehavi, E. and Wolf, J.K. (1987), 'On the performance evaluation of trellis codes', *IEEE Trans. Inf. Theory*, Vol. IT-33, No. 2, March, pp. 196-201.